

Probabilidad y estadística

para ingeniería y ciencias



NOVENA EDICIÓN

Walpole • Myers • Myers

Probabilidad y estadística
para ingeniería y ciencias

Segunda edición

Rodrigo C. Wazpale

Probabilidad y estadística para ingeniería y ciencias

MFN: 0000011464

579.5
.U35

Pro.
2012

← PROBABILIDAD →
← ESTADÍSTICA INFERENCIAL →



UNIVERSIDAD TECNICA DEL NORTE BIBLIOTECA	
Via de adquisición:	Compra
Documento N°:	082-PG-2013
Fecha:	22-11-2013
Valor unitario:	43,78
Código de Barras:	035392
MTCO	

Probabilidad y estadística para ingeniería y ciencias

Novena edición

Ronald E. Walpole
Roanoke College

Raymond H. Myers
Virginia Tech

Sharon L. Myers
Radford University

Keying Ye
University of Texas at San Antonio

Traducción
Leticia Esther Pineda Ayala
Traductora especialista en estadística

Revisión técnica
Roberto Hernández Ramírez
Departamento de Física y Matemáticas
División de Ingeniería y Tecnologías
Universidad de Monterrey

Linda Margarita Medina Herrera
Departamento de Física y Matemáticas
Escuela de Diseño, Ingeniería y Arquitectura
Instituto Tecnológico y de Estudios Superiores de Monterrey,
Campus Ciudad de México

PEARSON

Datos de catalogación bibliográfica

RONALD E. WALPOLE, RAYMOND H. MYERS,
SHARON L. MYERS Y KEYING YE

Probabilidad y estadística para ingeniería y ciencias
Novena edición

PEARSON EDUCACIÓN, México, 2012

ISBN: 978-607-32-1417-9

Área: Ingeniería

Formato: 18.5 × 23.5 cm

Páginas: 816

Authorized translation from the English language edition, entitled *PROBABILITY & STATISTICS FOR ENGINEERS & SCIENTISTS 9th Edition*, by *RONALD E. WALPOLE, RAYMOND H. MYERS, SHARON L. MYERS and KEYING YE*, published by Pearson Education, Inc., publishing as Pearson, Copyright © 2012. All rights reserved.
ISBN 9780321629111

Traducción autorizada de la edición en idioma inglés, titulada *PROBABILIDAD Y ESTADÍSTICA PARA INGENIERÍA Y CIENCIAS 9ª edición* por *RONALD E. WALPOLE, RAYMOND H. MYERS, SHARON L. MYERS y KEYING YE*, publicada por Pearson Education, Inc., publicada como Pearson, Copyright © 2012. Todos los derechos reservados.

Esta edición en español es la única autorizada.

Edición en español

Dirección Educación Superior: Mario Contreras

Editor sponsor:

Gabriela López Ballesteros

e-mail: gabriela.lopezballesteros@pearson.com

Editor de desarrollo:

Felipe Hernández Carrasco

Supervisor de Producción:

Juan José García Guzmán

Diseño de portada:

Dream Studio/Edgar Maldonado

Gerencia editorial

Educación Superior Latinoamérica: Marisa de Anta

NOVENA EDICIÓN, 2012

D.R. © 2012 por Pearson Educación de México, S.A. de C.V.

Atacomulco 500-5o. piso

Col. Industrial Atoto

53519, Naucalpan de Juárez, Estado de México

Cámara Nacional de la Industria Editorial Mexicana. Reg. núm. 1031.

Reservados todos los derechos. Ni la totalidad ni parte de esta publicación pueden reproducirse, registrarse o transmitirse, por un sistema de recuperación de información, en ninguna forma ni por ningún medio, sea electrónico, mecánico, fotoquímico, magnético o electroóptico, por fotocopia, grabación o cualquier otro, sin permiso previo por escrito del editor.

El préstamo, alquiler o cualquier otra forma de cesión de uso de este ejemplar requerirá también la autorización del editor o de sus representantes.

ISBN VERSIÓN IMPRESA: 978-607-32-1417-9

ISBN VERSIÓN E-BOOK: 978-607-32-1418-6

ISBN E-CHAPTER: 978-607-32-1419-3

Impreso en México. *Printed in Mexico.*

1 2 3 4 5 6 7 8 9 0 - 15 14 13 12

Esta obra se terminó de imprimir el mes de febrero de 2013
en los talleres de Editorial Progreso, S. A. de C. V.
Naranjo Núm. 248, Colonia Santa María la Ribera
Delegación Cuauhtémoc, C. P. 06400, México, D. F.

PEARSON

www.pearsonenespañol.com

AGRADECIMIENTOS

Pearson agradece a los profesores usuarios de esta obra y a los centros de estudio por su apoyo y retroalimentación, elementos fundamentales para esta nueva edición de Probabilidad y estadística para ingeniería y ciencias.

COLOMBIA

Escuela Colombiana de Ingeniería
Departamento de Matemáticas
Susana Rondón Troncoso

Pontificia Universidad Javeriana
Cali

Departamento de Ciencias
Naturales y Matemáticas
Daniel Enrique González Gómez
María del Pilar Marín Gaviria
Sandra Milena Ramírez Buelvas

Universidad Católica de Colombia
Departamento de Ciencias Básicas
Queeny Madueño Pinto

Universidad de La Salle
Departamento de Ciencias Básicas
Maribel Méndez Cortés
Martha Tatiana Jiménez Valderrama
Milton Armando Reyes Villamil
Myrian Elena Vergara Morales

COSTA RICA

Instituto Tecnológico de Costa Rica
Escuela de Ingeniería en
Producción Industrial
Ivannia Hasbum Fernández

Universidad de Costa Rica
Escuela de Estadística
Facultad de Ciencias Económicas
Ana Teresa Garita Salas

MÉXICO

Estado de México

Facultad de Estudios Superiores
Cuautitlán C-4
Armando Aguilar Márquez
Fermín Cervantes Martínez
Héctor Coss Garduño
Juan Carlos Axotla García
Miguel de Nazareth Pineda Becerril
Vicente Vázquez Juárez

Tecnológico de Estudios Superiores
de Coacalco
María de la Luz Dávila Flores
Martha Nieto López
Héctor Feliciano Martínez Osorio
Jeanette López Alanís

Tecnológico de Estudios Superiores
de Ecatepec
Héctor Rodríguez Carmona
Ángel Hernández Estrada
Daniel Jaimes Serrano
Ramón Jordán Rocha

Jalisco

Universidad de Guadalajara
Centro Universitario de Ciencias
Exactas e Ingenierías (CUCEI)
Departamento de Matemáticas
Agustín Rodríguez Martínez
Carlos Florentino Melgodo Cañedo
Cecilia Garibay López
Dalmiro García Nava

Deliazar Pantoja Espinoza
Gloria Arroyo Cervantes
Javier Nava Gómez
Jorge Luis Rodríguez Gutiérrez
José Ángel Partida Ibarra
José de Jesús Bernal Casillas
José de Jesús Cabrera Chavarría
José de Jesús Rivera Prado
José Solís Rodríguez
Julieta Carrasco García
Laura Esther Cortés Navarro
Lizbeth Díaz Caldera
Maribel Sierra Fuentes
Mario Alberto Prado Alonso
Oswaldo Camacho Castillo
Rosalfá Buenrostro Arceo
Samuel Rosalío Cuevas

Universidad del Valle de México,
Zapopan
Departamento de Ingeniería
Abel Vázquez Pérez
Irene Isabel Navarro González
Jorge Eduardo Aguilar Rosas
Miguel Arturo Barreiro González

Sinaloa

Instituto Tecnológico de Culiacán
Ciencias Básicas
Cecilia Norzagaray Gámez

Instituto Tecnológico de Los Mochis
Ciencias Básicas
Jesús Alberto Báez Torres

Contenido

Prefacio	XV
1 Introducción a la estadística y al análisis de datos.....	1
1.1 Panorama general: inferencia estadística, muestras, poblaciones y el papel de la probabilidad	1
1.2 Procedimientos de muestreo; recolección de los datos.....	7
1.3 Medidas de localización: la media y la mediana de una muestra	11
Ejercicios.....	13
1.4 Medidas de variabilidad.....	14
Ejercicios.....	17
1.5 Datos discretos y continuos	17
1.6 Modelado estadístico, inspección científica y diagnósticos gráficos.....	18
1.7 Tipos generales de estudios estadísticos: diseño experimental, estudio observacional y estudio retrospectivo	27
Ejercicios.....	30
2 Probabilidad	35
2.1 Espacio muestral.....	35
2.2 Eventos.....	38
Ejercicios.....	42
2.3 Conteo de puntos muestrales	44
Ejercicios.....	51
2.4 Probabilidad de un evento.....	52
2.5 Reglas aditivas	56
Ejercicios.....	59
2.6 Probabilidad condicional, independencia y regla del producto	62
Ejercicios.....	69
2.7 Regla de Bayes.....	72
Ejercicios.....	76
Ejercicios de repaso	77

2.8	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	79
3	Variables aleatorias y distribuciones de probabilidad	81
3.1	Concepto de variable aleatoria.....	81
3.2	Distribuciones discretas de probabilidad	84
3.3	Distribuciones de probabilidad continua	87
	Ejercicios.....	91
3.4	Distribuciones de probabilidad conjunta	94
	Ejercicios.....	104
	Ejercicios de repaso	107
3.5	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	109
4	Esperanza matemática.....	111
4.1	Media de una variable aleatoria	111
	Ejercicios.....	117
4.2	Varianza y covarianza de variables aleatorias.....	119
	Ejercicios.....	127
4.3	Medias y varianzas de combinaciones lineales de variables aleatorias	128
4.4	Teorema de Chebyshev	135
	Ejercicios.....	137
	Ejercicios de repaso	139
4.5	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	142
5	Algunas distribuciones de probabilidad discreta	143
5.1	Introducción y motivación	143
5.2	Distribuciones binomial y multinomial	143
	Ejercicios.....	150
5.3	Distribución hipergeométrica.....	152
	Ejercicios.....	157
5.4	Distribuciones binomial negativa y geométrica.....	158
5.5	Distribución de Poisson y proceso de Poisson.....	161
	Ejercicios.....	164
	Ejercicios de repaso	166
5.6	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	169

6	Algunas distribuciones continuas de probabilidad	171
6.1	Distribución uniforme continua	171
6.2	Distribución normal	172
6.3	Áreas bajo la curva normal	176
6.4	Aplicaciones de la distribución normal	182
	Ejercicios.....	185
6.5	Aproximación normal a la binomial	187
	Ejercicios.....	193
6.6	Distribución gamma y distribución exponencial	194
6.7	Distribución chi cuadrada	200
6.8	Distribución beta.....	201
6.9	Distribución logarítmica normal.....	201
6.10	Distribución de Weibull (opcional).....	203
	Ejercicios.....	206
	Ejercicios de repaso	207
6.11	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	209
7	Funciones de variables aleatorias (opcional).....	211
7.1	Introducción.....	211
7.2	Transformaciones de variables	211
7.3	Momentos y funciones generadoras de momentos.....	218
	Ejercicios.....	222
8	Distribuciones de muestreo fundamentales y descripciones de datos.....	225
8.1	Muestreo aleatorio	225
8.2	Algunos estadísticos importantes	227
	Ejercicios.....	230
8.3	Distribuciones muestrales	232
8.4	Distribución muestral de medias y el teorema del límite central.....	233
	Ejercicios.....	241
8.5	Distribución muestral de S^2	243
8.6	Distribución t	246
8.7	Distribución F	251
8.8	Gráficas de cuantiles y de probabilidad	254
	Ejercicios.....	259
	Ejercicios de repaso	260

8.9	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	262
9	Problemas de estimación de una y dos muestras	265
9.1	Introducción	265
9.2	Inferencia estadística	265
9.3	Métodos de estimación clásicos.....	266
9.4	Una sola muestra: estimación de la media.....	269
9.5	Error estándar de una estimación puntual.....	276
9.6	Intervalos de predicción.....	277
9.7	Límites de tolerancia.....	280
	Ejercicios.....	282
9.8	Dos muestras: estimación de la diferencia entre dos medias.....	285
9.9	Observaciones pareadas	291
	Ejercicios.....	294
9.10	Una sola muestra: estimación de una proporción	296
9.11	Dos muestras: estimación de la diferencia entre dos proporciones	300
	Ejercicios	302
9.12	Una sola muestra: estimación de la varianza	303
9.13	Dos muestras: estimación de la proporción de dos varianzas.....	305
	Ejercicios.....	307
9.14	Estimación de la máxima verosimilitud (opcional)	307
	Ejercicios.....	312
	Ejercicios de repaso	313
9.15	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	316
10	Pruebas de hipótesis de una y dos muestras.....	319
10.1	Hipótesis estadísticas: conceptos generales.....	319
10.2	Prueba de una hipótesis estadística.....	321
10.3	Uso de valores P para la toma de decisiones en la prueba de hipótesis	331
	Ejercicios.....	334
10.4	Una sola muestra: pruebas respecto a una sola media.....	336
10.5	Dos muestras: pruebas sobre dos medias.....	342
10.6	Elección del tamaño de la muestra para la prueba de medias	349
10.7	Métodos gráficos para comparar medias	354
	Ejercicios.....	356
10.8	Una muestra: prueba sobre una sola proporción.....	361
10.9	Dos muestras: pruebas sobre dos proporciones	363
	Ejercicios.....	365
10.10	Pruebas de una y dos muestras referentes a varianzas.....	366
	Ejercicios.....	369

10.11	Prueba de la bondad de ajuste.....	371
10.12	Prueba de independencia (datos categóricos).....	374
10.13	Prueba de homogeneidad.....	376
10.14	Estudio de caso de dos muestras.....	380
	Ejercicios.....	382
	Ejercicios de repaso	384
10.15	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	387
11	Regresión lineal simple y correlación.....	389
11.1	Introducción a la regresión lineal.....	389
11.2	El modelo de regresión lineal simple (RLS).....	390
11.3	Mínimos-cuadrados y el modelo ajustado	394
	Ejercicios.....	398
11.4	Propiedades de los estimadores de mínimos cuadrados	400
11.5	Inferencias sobre los coeficientes de regresión.....	403
11.6	Predicción	408
	Ejercicios.....	411
11.7	Selección de un modelo de regresión	414
11.8	El método del análisis de varianza.....	414
11.9	Prueba para la linealidad de la regresión: datos con observaciones repetidas.....	416
	Ejercicios.....	421
11.10	Gráficas de datos y transformaciones	424
11.11	Estudio de caso de regresión lineal simple	428
11.12	Correlación	430
	Ejercicios.....	435
	Ejercicios de repaso	436
11.13	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	442
12	Regresión lineal múltiple y ciertos modelos de regresión no lineal	443
12.1	Introducción.....	443
12.2	Estimación de los coeficientes	444
12.3	Modelo de regresión lineal en el que se utilizan matrices	447
	Ejercicios.....	450
12.4	Propiedades de los estimadores de mínimos cuadrados	453
12.5	Inferencias en la regresión lineal múltiple	455
	Ejercicios.....	461
12.6	Selección de un modelo ajustado mediante la prueba de hipótesis	462

12.7	Caso especial de ortogonalidad (opcional)	467
	Ejercicios.....	471
12.8	Variables categóricas o indicadoras	472
	Ejercicios.....	476
12.9	Métodos secuenciales para la selección del modelo.....	476
12.10	Estudio de los residuales y violación de las suposiciones (verificación del modelo).....	482
12.11	Validación cruzada, C_p , y otros criterios para la selección del modelo	487
	Ejercicios.....	494
12.12	Modelos especiales no lineales para condiciones no ideales.....	496
	Ejercicios.....	500
	Ejercicios de repaso	501
12.13	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	506
13	Experimentos con un solo factor: generales	507
13.1	Técnica del análisis de varianza.....	507
13.2	La estrategia del diseño de experimentos	508
13.3	Análisis de varianza de un factor: diseño completamente aleatorizado (ANOVA de un factor).....	509
13.4	Pruebas de la igualdad de varias varianzas	516
	Ejercicios.....	518
13.5	Comparaciones de un grado de libertad.....	520
13.6	Comparaciones múltiples.....	523
	Ejercicios.....	529
13.7	Comparación de un conjunto de tratamientos en bloques	532
13.8	Diseños de bloques completos aleatorizados.....	533
13.9	Métodos gráficos y verificación del modelo	540
13.10	Transformaciones de datos en el análisis de varianza	543
	Ejercicios.....	545
13.11	Modelos de efectos aleatorios.....	547
13.12	Estudio de caso	551
	Ejercicios.....	553
	Ejercicios de repaso	555
13.13	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	559
14	Experimentos factoriales (dos o más factores)	561
14.1	Introducción.....	561
14.2	Interacción en el experimento de dos factores.....	562
14.3	Análisis de varianza de dos factores	565
	Ejercicios.....	575

14.4	Experimentos de tres factores.....	579
	Ejercicios.....	586
14.5	Experimentos factoriales para efectos aleatorios y modelos mixtos	588
	Ejercicios.....	592
	Ejercicios de repaso	594
14.6	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	596
15	Experimentos factoriales 2^k y fracciones	597
15.1	Introducción.....	597
15.2	El factorial 2^k : cálculo de efectos y análisis de varianza	598
15.3	Experimento factorial 2^t sin réplicas	604
	Ejercicios.....	609
15.4	Experimentos factoriales en un ajuste de regresión.....	612
15.5	El diseño ortogonal.....	617
	Ejercicios.....	625
15.6	Experimentos factoriales fraccionarios.....	626
15.7	Análisis de experimentos factoriales fraccionados.....	632
	Ejercicios.....	634
15.8	Diseños de fracciones superiores y de filtrado	636
15.9	Construcción de diseños de resolución III y IV, con 8, 16 y 32 puntos de diseño.....	637
15.10	Otros diseños de resolución III de dos niveles; los diseños de Plackett-Burman.....	638
15.11	Introducción a la metodología de superficie de respuesta	639
15.12	Diseño robusto de parámetros.....	643
	Ejercicios.....	652
	Ejercicios de repaso	653
15.13	Posibles riesgos y errores conceptuales; relación con el material de otros capítulos	654
16	Estadística no paramétrica.....	655
16.1	Pruebas no paramétricas	655
16.2	Prueba de rango con signo.....	660
	Ejercicios.....	663
16.3	Prueba de la suma de rangos de Wilcoxon	665
16.4	Prueba de Kruskal-Wallis	668
	Ejercicios.....	670
16.5	Pruebas de rachas.....	671
16.6	Límites de tolerancia.....	674

16.7	Coefficiente de correlación de rango	674
	Ejercicios.....	677
	Ejercicios de repaso	679
17	Control estadístico de la calidad.....	681
17.1	Introducción.....	681
17.2	Naturaleza de los límites de control.....	683
17.3	Objetivos de la gráfica de control	683
17.4	Gráficas de control para variables.....	684
17.5	Gráficas de control para atributos	697
17.6	Gráficas de control de cusum.....	705
	Ejercicios de repaso	706
18	Estadística bayesiana	709
18.1	Conceptos bayesianos	709
18.2	Inferencias bayesianas	710
18.3	Estimados bayesianos mediante la teoría de decisión	717
	Ejercicios.....	718
	Bibliografía	721
	Apéndice A: Tablas y demostraciones estadísticas.....	725
	Apéndice B: Respuestas a los ejercicios impares (no de repaso)	769
	Índice.....	785

Prefacio

Enfoque general y nivel matemático

Al elaborar la novena edición, nuestro interés principal no fue tan sólo agregar material nuevo sino brindar claridad y mejor comprensión. Este objetivo se logró en parte al incluir material nuevo al final de los capítulos, lo cual permite que se relacionen mejor. Con cierto afecto llamamos “contratiempos” a los comentarios que aparecen al final de los capítulos, pues son muy útiles para que los estudiantes recuerden la idea general y la forma en que cada capítulo se ajusta a esa imagen; así como para que entiendan las limitaciones y los problemas que resultarían por el uso inadecuado de los procedimientos. Los proyectos para la clase favorecen una mayor comprensión de cómo se utiliza la estadística en el mundo real, por lo que añadimos algunos proyectos en varios capítulos. Tales proyectos brindan a los estudiantes la oportunidad de trabajar solos o en equipo, y de reunir sus propios datos experimentales para realizar inferencias. En algunos casos, el trabajo implica un problema cuya solución ejemplifica el significado de un concepto, o bien, favorece la comprensión empírica de un resultado estadístico importante. Se ampliaron algunos de los ejemplos anteriores y se introdujeron algunos nuevos para crear “estudios de caso”, los cuales incluyen un comentario para aclarar al estudiante un concepto estadístico en el contexto de una situación práctica.

En esta edición seguimos haciendo énfasis en el equilibrio entre la teoría y las aplicaciones. Utilizamos el cálculo y otros tipos de conceptos matemáticos, por ejemplo, de álgebra lineal, casi al mismo nivel que en ediciones anteriores. Las herramientas analíticas para la estadística se cubren de mejor manera utilizando el cálculo en los casos donde el análisis se centra en las reglas de los conceptos de probabilidad. En los capítulos 2 a 10 se destacan las distribuciones de probabilidad y la inferencia estadística. En los capítulos 11 a 15, en los cuales se estudian la regresión lineal y el análisis de varianza, se aplica un poco de álgebra lineal y matrices. Los estudiantes que utilizan este libro deben haber cursado el equivalente a un semestre de cálculo diferencial e integral. El álgebra lineal es útil aunque no indispensable, siempre y cuando el instructor no cubra la sección sobre regresión lineal múltiple del capítulo 12 utilizando álgebra de matrices. Al igual que en las ediciones anteriores, y con la finalidad de desafiar al estudiante, muchos ejercicios se refieren a aplicaciones científicas y de ingeniería a la vida real. Todos los conjuntos de datos asociados con los ejercicios están disponibles para descargar del sitio web <http://www.pearsonespañol.com/walpole>.

Resumen de los cambios en la novena edición

- Para brindar una mayor comprensión del uso de la estadística en el mundo real, en varios capítulos se agregaron proyectos para la clase. Los estudiantes tienen que generar o reunir sus propios datos experimentales y realizar inferencias a partir de ellos.
- Se agregaron más estudios de caso y otros se ampliaron para ayudar a los usuarios a comprender los métodos estadísticos que se presentan en el contexto de una situación real. Por ejemplo, la interpretación de los límites de confianza, los límites de predicción y los límites de tolerancia se exponen utilizando situaciones de la vida real.
- Se agregaron "contratiempos" al final de algunos capítulos y en otros se ampliaron los que ya se incluían. El objetivo de dichos comentarios es presentar cada capítulo en el contexto de la idea general y analizar la forma en que los capítulos se relacionan entre sí. Otro objetivo es advertir acerca del uso inadecuado de las técnicas estadísticas examinadas en el capítulo.
- El capítulo 1 se mejoró y ahora incluye más estadísticos de una sola cifra y técnicas gráficas. También se incluyó nuevo material fundamental sobre muestreo y diseño experimental.
- Los ejemplos que se agregaron en el capítulo 8 sobre las distribuciones de muestreo tienen la finalidad de motivar a los estudiantes a realizar las pruebas de hipótesis y de los valores P . Esto los prepara para el material más avanzado sobre los temas que se presentan en el capítulo 10.
- El capítulo 12 contiene más información sobre el efecto que tiene una sola variable de regresión en un modelo que presenta una gran colinealidad con otras variables.
- El capítulo 15 ahora introduce material sobre el importante tema de la metodología de superficie de respuesta (MSR). El uso de las variables del ruido en la MSR permite ejemplificar los modelos de la media y la varianza (superficie de respuesta doble).
- En el capítulo 15 se introduce el diseño compuesto central.
- El capítulo 18 incluye más ejemplos y un mejor análisis de cómo se utilizan los métodos bayesianos para la toma de decisiones estadísticas.

Contenido y planeación del curso

Este libro está diseñado para un curso de uno o dos semestres. Un plan razonable para el curso de un semestre podría incluir los capítulos 1 a 10, lo cual daría como resultado un programa que concluye con los fundamentos de la estimación y la prueba de hipótesis. Los profesores que desean que los estudiantes aprendan la regresión lineal simple podrían incluir una parte del capítulo 11. Para quienes desean incluir el análisis de varianza en vez de la regresión, el curso de un semestre podría incluir el capítulo 13 en vez de los capítulos 11 y 12. El capítulo 13 trata el tema del análisis de varianza de un factor. Otra opción consiste en eliminar partes de los capítulos 5 o 6, así como el capítulo 7. Al hacer esto se omitirían las distribuciones discretas o continuas, mismas que incluyen la binomial negativa, la geométrica, la gamma, la de Weibull, la beta y la logarítmica normal. Otros contenidos que se podrían omitir en un programa de un semestre son la estimación de máxima verosimilitud, la predicción y los límites de tolerancia del

capítulo 9. El programa para un semestre suele ser flexible, dependiendo del interés que el profesor tenga en la regresión, el análisis de varianza, el diseño experimental y los métodos de superficie de respuesta (capítulo 15). Existen varias distribuciones discretas y continuas (capítulos 5 y 6) que tienen aplicaciones en diversas áreas de la ingeniería y las ciencias.

Los capítulos 11 a 18 incluyen una gran cantidad de material que se podría agregar al segundo semestre, en caso de que se eligiera un curso de dos semestres. El material sobre la regresión lineal simple y múltiple se estudia en los capítulos 11 y 12, respectivamente. El capítulo 12 puede ser muy flexible. La regresión lineal múltiple incluye “temas especiales”, como variables categóricas o indicadoras, métodos secuenciales para la selección de modelos, por ejemplo, la regresión por etapas, el estudio de residuales para la detección de violaciones de supuestos, la validación cruzada y el uso de los estadísticos PRESS, así como el de C_p y la regresión logística. Se hace hincapié en el uso de regresores ortogonales, un precursor del diseño experimental en el capítulo 15. Los capítulos 13 y 14 ofrecen hasta cierto grado material abundante sobre el análisis de varianza (ANOVA), con modelos fijos, aleatorios y mixtos. En el capítulo 15 se destaca la aplicación de los diseños con dos niveles en el contexto de los experimentos factoriales fraccionarios y completos (2^k). También se ejemplifican los diseños especiales de selección. En el capítulo 15 se incluye asimismo una nueva sección sobre la metodología de superficie de respuesta (MSR), para ejemplificar el uso del diseño experimental con la finalidad de encontrar condiciones óptimas de proceso. Se analiza el ajuste de un modelo de segundo orden utilizando un diseño complejo central. La MSR se amplía para abarcar el análisis de problemas sobre el diseño de un parámetro robusto. Las variables de ruido se utilizan para ajustar modelos dobles de superficie de respuesta. Los capítulos 16, 17 y 18 incluyen una cantidad moderada de material sobre estadística no paramétrica, control de calidad e inferencia bayesiana.

El capítulo 1 es un bosquejo de la inferencia estadística, presentada a un nivel matemático sencillo, pero de manera más amplia que en la octava edición con el propósito de examinar más detalladamente los estadísticos de una sola cifra y las técnicas gráficas. Este capítulo está diseñado para brindar a los estudiantes una presentación preliminar de los conceptos fundamentales que les permitirán entender los detalles posteriores de mayor complejidad. Se presentan conceptos clave sobre muestreo, recolección de datos y diseño experimental, así como los aspectos rudimentarios de las herramientas gráficas y la información que se obtiene a partir de un conjunto de datos. También se agregaron las gráficas de tallo y hojas, y las de caja y bigotes. Las gráficas están mejor organizadas y etiquetadas. El análisis de la incertidumbre y la variación en un sistema se ilustra de forma detallada. Se incluyen ejemplos de cómo clasificar las características importantes de un sistema o proceso científico, y esas ideas se ilustran en ambientes prácticos, como procesos de manufactura, estudios biomédicos, y estudios de sistemas biológicos y científicos de otros tipos. Se efectúa una comparación entre el uso de los datos discretos y continuos; también se hace un mayor énfasis en el uso de modelos y de la información con respecto a los modelos estadísticos que se logran obtener mediante las herramientas gráficas.

En los capítulos 2, 3 y 4 se estudian los conceptos básicos de probabilidad, así como las variables aleatorias discretas y continuas. Los capítulos 5 y 6 se enfocan en las distribuciones discretas y continuas específicas, así como en las relaciones que existen entre ellas. En estos capítulos también se destacan ejemplos de aplicaciones de las distribuciones en estudios reales científicos y de ingeniería. Los estudios de caso, los ejemplos y una gran cantidad de ejercicios permiten a los estudiantes practicar el uso de tales distribuciones. Los proyectos permiten la aplicación práctica de estas distribuciones en la vida

real mediante el trabajo en equipo. El capítulo 7 es el más teórico del libro; en él se expone la transformación de variables aleatorias, y podría ser que no se utilice a menos que el instructor desee impartir un curso relativamente teórico. El capítulo 8 contiene material gráfico, el cual amplía el conjunto básico de herramientas gráficas presentadas y ejemplificadas en el capítulo 1. Aquí se analizan las gráficas de probabilidad y se ilustran con ejemplos. El muy importante concepto de las distribuciones de muestreo se presenta de forma detallada, y se proporcionan ejemplos que incluyen el teorema del límite central y la distribución de una varianza muestral en una situación de muestreo independiente y normal. También se presentan las distribuciones t y F para motivar a los estudiantes a utilizarlas en los capítulos posteriores. El nuevo material del capítulo 8 ayuda a los estudiantes a conocer la importancia de la prueba de hipótesis mediante la presentación del concepto del valor P .

El capítulo 9 contiene material sobre la estimación puntual y de intervalos de una muestra y dos muestras. Un análisis detallado y con ejemplos destaca las diferencias entre los tipos de intervalos (intervalos de confianza, intervalos de predicción e intervalos de tolerancia). Un estudio de caso ilustra los tres tipos de intervalos estadísticos en el contexto de una situación de manufactura. Este estudio de caso destaca las diferencias entre los intervalos, sus fuentes y los supuestos en que se basan, así como cuáles son los intervalos que requieren diferentes tipos de estudios o preguntas. Se añadió un método de aproximación para las inferencias sobre una proporción. El capítulo 10 inicia con una presentación básica sobre el significado práctico de la prueba de hipótesis, con un énfasis en conceptos fundamentales como la hipótesis nula y la alternativa, el papel que desempeñan la probabilidad y el valor P , así como la potencia de una prueba. Después, se presentan ejemplos de pruebas sobre una o dos muestras en condiciones estándar. También se describe la prueba t de dos muestras con observaciones en pares (apareadas). Un estudio de caso ayuda a los estudiantes a entender el verdadero significado de una interacción de factores, así como los problemas que en ocasiones surgen cuando existen interacciones entre tratamientos y unidades experimentales. Al final del capítulo 10 se incluye una sección muy importante que relaciona los capítulos 9 y 10 (estimación y prueba de hipótesis) con los capítulos 11 a 16, donde se destaca el modelamiento estadístico. Es importante que el estudiante esté consciente de la fuerte relación entre los capítulos mencionados.

Los capítulos 11 y 12 incluyen material sobre la regresión lineal simple y múltiple, respectivamente. En esta edición ponemos mucho más atención en el efecto que tiene la colinealidad entre las variables de regresión. Se presenta una situación que muestra cómo el papel que desempeña una sola variable de regresión depende en gran parte de cuáles son los regresores que la acompañan en el modelo. Después se revisan los procedimientos secuenciales para la selección del modelo (hacia adelante, hacia atrás, por etapas, etcétera) con respecto a este concepto, así como los fundamentos para utilizar ciertos tipos de valores P con tales procedimientos. En el capítulo 12 se estudia material sobre los modelos no lineales con una presentación especial de la regresión logística, la cual tiene aplicaciones en ingeniería y en las ciencias biológicas. El material sobre la regresión múltiple es muy extenso, de manera que, como antes se expuso, plantea una gran flexibilidad. Al final del capítulo 12 se incluye un comentario que lo relaciona con los capítulos 14 y 15. Se agregaron varios elementos para fomentar la comprensión del material en general. Por ejemplo, al final del capítulo se describen algunas dificultades y problemas que podrían surgir. Se indica que existen tipos de respuestas que ocurren de forma natural en la práctica, por ejemplo, respuestas de proporciones, de conteo y muchas otras, con las cuales no se debe utilizar la regresión estándar de mínimos cuadrados

debido a que los supuestos de normalidad no se cumplen, y transgredirlos causaría errores muy graves. Se sugiere utilizar la transformación de datos para reducir el problema en algunos casos. Nuevamente, los capítulos 13 y 14 sobre el tema del análisis de varianza tienen cierta flexibilidad. En el capítulo 13 se estudia el ANOVA de un factor en el contexto de un diseño completamente aleatorio. Algunos temas complementarios incluyen las pruebas sobre las varianzas y las comparaciones múltiples. Se destacan las comparaciones de tratamientos en bloque, junto con el tema de los bloques completos aleatorizados. Los métodos gráficos se extendieron al ANOVA para ayudar al estudiante a complementar la inferencia formal con una inferencia pictórica que facilita la presentación del material a los científicos y a los ingenieros. Se incluye un nuevo proyecto donde los estudiantes incorporan la aleatoriedad adecuada a cada plan, y se utilizan técnicas gráficas y valores P en el informe de los resultados. En el capítulo 14 se amplía el material del capítulo 13 para ajustar dos o más factores dentro de una estructura factorial. La presentación del ANOVA en el capítulo 14 incluye la creación de modelos aleatorios y de efectos fijos. En el capítulo 15 se estudia material relacionado con los diseños factoriales 2^k ; los ejemplos y los estudios de caso plantean el uso de diseños de selección y fracciones especiales de orden superior del factorial 2^k . Dos elementos nuevos y especiales son la metodología de superficie de respuesta (MSR) y el diseño de parámetros robustos. Son temas que se relacionan en un estudio de caso que describe e ilustra un diseño doble de superficie de respuesta, así como un análisis que incluye el uso de superficies de respuesta de la media y la varianza de procesos.

Programa de cómputo

Los estudios de caso, que inician en el capítulo 8, muestran impresiones de listas de resultados por computadora y material gráfico generado con los programas SAS y MINITAB. El hecho de incluir los cálculos por computadora refleja nuestra idea de que los estudiantes deben contar con la experiencia de leer e interpretar impresiones de listas de resultados y gráficas por computadora, incluso si el software que se utiliza en el libro no coincide con el que utiliza el profesor. La exposición a más de un tipo de programas aumentaría la experiencia de los estudiantes. No hay razones para creer que el programa utilizado en el curso coincidirá con el que el estudiante tendrá que utilizar en la práctica después de graduarse. Cuando sea pertinente, los ejemplos y los estudios de caso en el libro se complementarán con diversos tipos de gráficas residuales, cuantilares, de probabilidad normal y de otros tipos. Tales gráficas se incluyen especialmente en los capítulos 11 a 15.

Complementos

Manual de soluciones para el instructor. Este recurso contiene respuestas a todos los ejercicios del libro y se puede descargar del Centro de Recursos para Profesor de Pearson.

Diapositivas de PowerPoint® ISBN-10: 0-321-73731-8; ISBN-13: 978-0-321-73731-1. Las diapositivas incluyen la mayoría de las figuras y las tablas del libro; se pueden descargar del Centro de Recursos para el Profesor de Pearson.

Reconocimientos

Estamos en deuda con los colegas que revisaron las anteriores ediciones de este libro y que nos dieron muchas sugerencias útiles para esta edición. Ellos son David Groggel, de *Miami University*; Lance Hemlow, de *Raritan Valley Community College*; Ying Ji, de *University of Texas at San Antonio*; Thomas Kline, de *University of Northern Iowa*; Sheila Lawrence, de *Rutgers University*; Luis Moreno, de *Broome County Community College*; Donald Waldman, de *University of Colorado-Boulder* y Marlene Will, de *Spalding University*. También queremos agradecer a Delray Schulz, de *Millersville University*, Roxane Burrows, de *Hocking College* y Frank Chmely por asegurarse de la exactitud de este libro.

Nos gustaría agradecer a la editorial y a los servicios de producción suministrados por muchas personas de Pearson/Prentice Hall, sobre todo a Deirdre Lynch, la editora en jefe, a Christopher Cummings, el editor de adquisiciones, a Christine O'Brien, la editora de contenido ejecutivo, a Tracy Patruno, la editora de producción y a Sally Lifland, la editora de producción. Apreciamos los comentarios y sugerencias útiles de Gail Magin, la correctora de estilo. También estamos en deuda con el Centro de Asesoría Estadística de Virginia Tech, que fue nuestra fuente de muchos conjuntos reales de datos.

R.H.M.
S.L.M.
K.Y.

CAPÍTULO 1

Introducción a la estadística y al análisis de datos

1.1 Panorama general: inferencia estadística, muestras, poblaciones y el papel de la probabilidad

Desde inicios de la década de los ochenta del siglo pasado y hasta lo que ha transcurrido del siglo XXI la industria estadounidense ha puesto una enorme atención en el *mejoramiento de la calidad*. Se ha dicho y escrito mucho acerca del “milagro industrial” en Japón, que comenzó a mediados del siglo XX. Los japoneses lograron el éxito en donde otras naciones fallaron, a saber, en la creación de un entorno que permita la manufactura de productos de alta calidad. Gran parte del éxito de los japoneses se atribuye al uso de *métodos estadísticos* y del pensamiento estadístico entre el personal gerencial.

Empleo de datos científicos

El uso de métodos estadísticos en la manufactura, el desarrollo de productos alimenticios, el software para computadoras, las fuentes de energía, los productos farmacéuticos y muchas otras áreas implican el acopio de información o **datos científicos**. Por supuesto que la obtención de datos no es algo nuevo, ya que se ha realizado por más de mil años. Los datos se han recabado, resumido, reportado y almacenado para su examen cuidadoso. Sin embargo, hay una diferencia profunda entre el acopio de información científica y la **estadística inferencial**. Esta última ha recibido atención legítima en décadas recientes.

La estadística inferencial generó un número enorme de “herramientas” de los métodos estadísticos que utilizan los profesionales de la estadística. Los métodos estadísticos se diseñan para contribuir al proceso de realizar juicios científicos frente a la **incertidumbre** y a la **variación**. Dentro del proceso de manufactura, la densidad de producto de un material específico no siempre será la misma. De hecho, si un proceso es discontinuo en vez de continuo, la densidad de material no sólo variará entre los lotes que salen de la línea de producción (variación de un lote a otro), sino también dentro de los propios lotes. Los métodos estadísticos se utilizan para analizar datos de procesos como el anterior; el objetivo de esto es tener una mejor orientación respecto de cuáles cambios se deben realizar en el proceso para mejorar su **calidad**. En este proceso la calidad bien podría

definirse en relación con su grado de acercamiento a un valor de densidad meta en armonía con *qué parte de las veces* se cumple este criterio de cercanía. A un ingeniero podría interesarle un instrumento específico que se utilice para medir el monóxido de azufre en estudios sobre la contaminación atmosférica. Si el ingeniero dudara respecto de la eficacia del instrumento, tendría que tomar en cuenta dos **fuentes de variación**. La primera es la variación en los valores del monóxido de azufre que se encuentran en el mismo lugar el mismo día. La segunda es la variación entre los valores observados y la cantidad **real** de monóxido de azufre que haya en el aire en ese momento. Si cualquiera de estas dos fuentes de variación es excesivamente grande (según algún estándar determinado por el ingeniero), quizá se necesite remplazar el instrumento. En un estudio biomédico de un nuevo fármaco que reduce la hipertensión, 85% de los pacientes experimentaron alivio; aunque por lo general se reconoce que el medicamento actual o el “viejo” alivia a 80% de los pacientes que sufren hipertensión crónica. Sin embargo, el nuevo fármaco es más caro de elaborar y podría tener algunos efectos colaterales. ¿Se debería adoptar el nuevo medicamento? Éste es un problema con el que las empresas farmacéuticas, junto con la FDA (Federal Drug Administration), se encuentran a menudo (a veces es mucho más complejo). De nuevo se debe tomar en cuenta las necesidades de variación. El valor del “85%” se basa en cierto número de pacientes seleccionados para el estudio. Tal vez si se repitiera el estudio con nuevos pacientes ¡el número observado de “éxitos” sería de 75%! Se trata de una variación natural de un estudio a otro que se debe tomar en cuenta en el proceso de toma de decisiones. Es evidente que tal variación es importante, ya que la variación de un paciente a otro es endémica al problema.

Variabilidad en los datos científicos

En los problemas analizados anteriormente los métodos estadísticos empleados tienen que ver con la variabilidad y en cada caso la variabilidad que se estudia se encuentra en datos científicos. Si la densidad del producto observada en el proceso fuera siempre la misma y siempre fuera la esperada, no habría necesidad de métodos estadísticos. Si el dispositivo para medir el monóxido de azufre siempre diera el mismo valor y éste fuera exacto (es decir, correcto), no se requeriría análisis estadístico. Si entre un paciente y otro no hubiera variabilidad inherente a la respuesta al medicamento (es decir, si el fármaco siempre causara alivio o nunca aliviara), la vida sería muy sencilla para los científicos de las empresas farmacéuticas y de la FDA, y los estadísticos no serían necesarios en el proceso de toma de decisiones. Los investigadores de la estadística han originado un gran número de métodos analíticos que permiten efectuar análisis de datos obtenidos de sistemas como los descritos anteriormente, lo cual refleja la verdadera naturaleza de la ciencia que conocemos como estadística inferencial, a saber, el uso de técnicas que, al permitimos obtener conclusiones (o inferencias) sobre el sistema científico, nos permiten ir más allá de sólo reportar datos. Los profesionales de la estadística usan leyes fundamentales de probabilidad e inferencia estadística para sacar conclusiones respecto de los sistemas científicos. La información se colecta en forma de **muestras** o conjuntos de **observaciones**. En el capítulo 2 se introduce el proceso de muestreo, el cual se continúa analizando a lo largo de todo el libro.

Las muestras se reúnen a partir de **poblaciones**, que son conjuntos de todos los individuos o elementos individuales de un tipo específico. A veces una población representa un sistema científico. Por ejemplo, un fabricante de tarjetas para computadora podría desear eliminar defectos. Un proceso de muestreo implicaría recolectar información de 50 tarjetas de computadora tomadas aleatoriamente durante el proceso. En este caso la población

sería representada por todas las tarjetas de computadora producidas por la empresa en un periodo específico. Si se lograra mejorar el proceso de producción de las tarjetas para computadora y se reuniera una segunda muestra de tarjetas, cualquier conclusión que se obtuviera respecto de la efectividad del cambio en el proceso debería extenderse a toda la población de tarjetas para computadora que se produzcan en el “proceso mejorado”. En un experimento con fármacos se toma una muestra de pacientes y a cada uno se le administra un medicamento específico para reducir la presión sanguínea. El interés se enfoca en obtener conclusiones sobre la población de quienes sufren hipertensión. A menudo, cuando la planeación ocupa un lugar importante en la agenda, es muy importante el acopio de datos científicos en forma sistemática. En ocasiones la planeación está, por necesidad, bastante limitada. Con frecuencia nos enfocamos en ciertas propiedades o características de los elementos u objetos de la población. Cada característica tiene importancia de ingeniería específica o, digamos, biológica para el “cliente”, el científico o el ingeniero que busca aprender algo acerca de la población. Por ejemplo, en uno de los casos anteriores la calidad del proceso se relacionaba con la densidad del producto al salir del proceso. Un(a) ingeniero(a) podría necesitar estudiar el efecto de las condiciones del proceso, la temperatura, la humedad, la cantidad de un ingrediente particular, etcétera. Con ese fin podría mover de manera sistemática estos **factores** a cualesquiera niveles que se sugieran, de acuerdo con cualquier prescripción o **diseño experimental** que se desee. Sin embargo, un científico silvicultor que está interesado en estudiar los factores que influyen en la densidad de la madera en cierta clase de árbol no necesariamente tiene que diseñar un experimento. Este caso quizá requiera un **estudio observacional**, en el cual los datos se acopian en el campo pero no es posible seleccionar de antemano los **niveles de los factores**. Ambos tipos de estudio se prestan a los métodos de la inferencia estadística. En el primero, la calidad de las inferencias dependerá de la planeación adecuada del experimento. En el segundo, el científico está a expensas de lo que pueda recopilar. Por ejemplo, si un agrónomo se interesara en estudiar el efecto de la lluvia sobre la producción de plantas sería lamentable que recopilara los datos durante una sequía.

Es bien conocida la importancia del pensamiento estadístico para los administradores y el uso de la inferencia estadística para el personal científico. Los investigadores obtienen mucho de los datos científicos. Los datos proveen conocimiento acerca del fenómeno científico. Los ingenieros de producto y de procesos aprenden más en sus esfuerzos fuera de línea para mejorar el proceso. También logran una comprensión valiosa al reunir datos de producción (supervisión en línea) sobre una base regular, lo cual les permite determinar las modificaciones que se requiere realizar para mantener el proceso en el nivel de calidad deseado.

En ocasiones un científico sólo desea obtener alguna clase de resumen de un conjunto de datos representados en la muestra. En otras palabras, no requiere estadística inferencial. En cambio, le sería útil un conjunto de estadísticos o la **estadística descriptiva**. Tales números ofrecen un sentido de la ubicación del centro de los datos, de la variabilidad en los datos y de la naturaleza general de la distribución de observaciones en la muestra. Aunque no se incorporen métodos estadísticos específicos que lleven a la **inferencia estadística**, se puede aprender mucho. A veces la estadística descriptiva va acompañada de gráficas. El software estadístico moderno permite el cálculo de **medias, medianas, desviaciones estándar** y otros estadísticos de una sola cifra, así como el desarrollo de gráficas que presenten una “huella digital” de la naturaleza de la muestra. En las secciones siguientes veremos definiciones e ilustraciones de los estadísticos y descripciones de recursos gráficos como histogramas, diagramas de tallo y hojas, diagramas de dispersión, gráficas de puntos y diagramas de caja.

El papel de la probabilidad

En los capítulos 2 a 6 de este libro se presentan los conceptos fundamentales de la probabilidad. Un estudio concienzudo de las bases de tales conceptos permitirá al lector comprender mejor la inferencia estadística. Sin algo de formalismo en teoría de probabilidad, el estudiante no podría apreciar la verdadera interpretación del análisis de datos a través de los métodos estadísticos modernos. Es muy natural estudiar probabilidad antes de estudiar inferencia estadística. Los elementos de probabilidad nos permiten cuantificar la fortaleza o “confianza” en nuestras conclusiones. En este sentido, los conceptos de probabilidad forman un componente significativo que complementa los métodos estadísticos y ayuda a evaluar la consistencia de la inferencia estadística. Por consiguiente, la disciplina de la probabilidad brinda la transición entre la estadística descriptiva y los métodos inferenciales. Los elementos de la probabilidad permiten expresar la conclusión en el lenguaje que requieren los científicos y los ingenieros. El ejemplo que sigue permite al lector comprender la noción de un valor- P , el cual a menudo proporciona el “fundamento” en la interpretación de los resultados a partir del uso de métodos estadísticos.

Ejemplo 1.1: Suponga que un ingeniero se encuentra con datos de un proceso de producción en el cual se muestrean 100 artículos y se obtienen 10 defectuosos. Se espera y se anticipa que ocasionalmente habrá artículos defectuosos. Obviamente estos 100 artículos representan la muestra. Sin embargo, se determina que, a largo plazo, la empresa sólo puede tolerar 5% de artículos defectuosos en el proceso. Ahora bien, los elementos de probabilidad permiten al ingeniero determinar qué tan concluyente es la información muestral respecto de la naturaleza del proceso. En este caso la **población** representa conceptualmente todos los artículos posibles en el proceso. Suponga que averiguamos que, *si el proceso es aceptable*, es decir, que su producción no excede un 5% de artículos defectuosos, hay una probabilidad de 0.0282 de obtener 10 o más artículos defectuosos en una muestra aleatoria de 100 artículos del proceso. Esta pequeña probabilidad sugiere que, en realidad, a largo plazo el proceso tiene un porcentaje de artículos defectuosos mayor al 5%. En otras palabras, en las condiciones de un proceso aceptable casi nunca se obtendría la información muestral que se obtuvo. Sin embargo, ¡se obtuvo! Por lo tanto, es evidente que la probabilidad de que se obtuviera sería mucho mayor si la tasa de artículos defectuosos del proceso fuera mucho mayor que 5%. ▮

A partir de este ejemplo se vuelve evidente que los elementos de probabilidad ayudan a traducir la información muestral en algo concluyente o no concluyente acerca del sistema científico. De hecho, lo aprendido probablemente constituya información inquietante para el ingeniero o administrador. Los métodos estadísticos (que examinaremos con más detalle en el capítulo 10) produjeron un valor- P de 0.0282. El resultado sugiere que es **muy probable que el proceso no sea aceptable**. En los capítulos siguientes se trata detenidamente el concepto de valor- P . El próximo ejemplo brinda una segunda ilustración.

Ejemplo 1.2: Con frecuencia, la naturaleza del estudio científico señalará el papel que desempeñan la probabilidad y el razonamiento deductivo en la inferencia estadística. El ejercicio 9.40 en la página 294 proporciona datos asociados con un estudio que se llevó a cabo en el Virginia Polytechnic Institute and State University acerca del desarrollo de una relación entre las raíces de los árboles y la acción de un hongo. Los minerales de los hongos se transfieren a los árboles, y los azúcares de los árboles a los hongos. Se plantaron dos muestras de 10 plántones de roble rojo norteamericano en un invernadero, una de ellas contenía

plantones tratados con nitrógeno y la otra plantones sin tratamiento. Todas las demás condiciones ambientales se mantuvieron constantes. Todos los plantones contenían el hongo *Pisolithus tinctorius*. En el capítulo 9 se incluyen más detalles. Los pesos en gramos de los tallos se registraron después de 140 días y los datos se presentan en la tabla 1.1.

Tabla 1.1: Conjunto de datos del ejemplo 1.2

Sin nitrógeno	Con nitrógeno
0.32	0.26
0.53	0.43
0.28	0.47
0.37	0.49
0.47	0.52
0.43	0.75
0.36	0.79
0.42	0.86
0.38	0.62
0.43	0.46

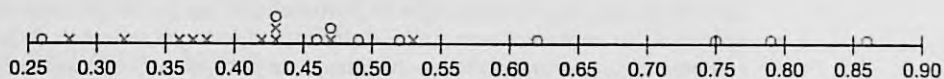


Figura 1.1: Gráfica de puntos de los datos de peso del tallo.

En este ejemplo hay dos muestras tomadas de dos poblaciones distintas. El objetivo del experimento es determinar si el uso del nitrógeno influye en el crecimiento de las raíces. Éste es un estudio comparativo (es decir, es un estudio en el que se busca comparar las dos poblaciones en cuanto a ciertas características importantes). Los datos se deben graficar como se indica en el diagrama de puntos de la figura 1.1. Los valores \circ representan los datos “con nitrógeno” y los valores \times los datos “sin nitrógeno”.

Observe que la apariencia general de los datos podría sugerir al lector que, en promedio, el uso del nitrógeno aumenta el peso del tallo. Cuatro observaciones con nitrógeno son considerablemente más grandes que cualquiera de las observaciones sin nitrógeno. La mayoría de las observaciones sin nitrógeno parece estar por debajo del centro de los datos. La apariencia del conjunto de datos parece indicar que el nitrógeno es efectivo. Pero, ¿cómo se cuantifica esto? ¿Cómo se puede resumir toda la evidencia visual aparente de manera que tenga algún significado? Como en el ejemplo anterior, se pueden utilizar los fundamentos de la probabilidad. Las conclusiones se resumen en una declaración de probabilidad o valor- P . Aquí no demostraremos la inferencia estadística que produce la probabilidad resumida. Igual que en el ejemplo 1.1, tales métodos se estudiarán en el capítulo 10. El problema gira alrededor de la “probabilidad de que datos como éstos se puedan observar”, *dado que el nitrógeno no tiene efecto*; en otras palabras, *dado que ambas muestras se generaron a partir de la misma población*. Suponga que esta probabilidad es pequeña, digamos de 0.03; un porcentaje que podría constituir suficiente evidencia de que el uso del nitrógeno en realidad influye en el peso promedio del tallo en los plantones de roble rojo (aparentemente lo aumenta). ▮

¿Cómo trabajan juntas la probabilidad y la inferencia estadística?

Es importante para el lector que comprenda claramente la diferencia entre la disciplina de la probabilidad, una ciencia por derecho propio, y la disciplina de la estadística inferencial. Como señalamos, el uso o la aplicación de conceptos de probabilidad permite interpretar la vida cotidiana a partir de los resultados de la inferencia estadística. En consecuencia, se afirma que la inferencia estadística emplea los conceptos de probabilidad. A partir de los dos ejemplos anteriores aprendimos que la información muestral está disponible para el analista y que, con la ayuda de métodos estadísticos y elementos de probabilidad, podemos obtener conclusiones acerca de alguna característica de la población (en el ejemplo 1.1 el proceso al parecer no es aceptable, y en el ejemplo 1.2 parece ser que el nitrógeno en verdad influye en el peso promedio de los tallos). Así, para un problema estadístico, **la muestra, junto con la estadística inferencial, nos permite obtener conclusiones acerca de la población, ya que la estadística inferencial utiliza ampliamente los elementos de probabilidad.** Tal razonamiento es *inductivo* por naturaleza. Ahora, cuando avancemos al capítulo 2 y los siguientes, el lector encontrará que, a diferencia de lo que hicimos en nuestros dos ejemplos actuales, no nos enfocaremos en resolver problemas estadísticos. En muchos de los ejemplos que estudiaremos no utilizaremos muestras. Lo que haremos será describir claramente una población con todas sus características conocidas. Las preguntas importantes se enfocarán en la naturaleza de los datos que hipotéticamente se podrían obtener a partir de la población. Entonces, podríamos afirmar que **los elementos de probabilidad nos permiten sacar conclusiones acerca de las características de los datos hipotéticos que se tomen de la población, con base en las características conocidas de la población.** Esta clase de razonamiento es *deductivo* por naturaleza. La figura 1.2 muestra la relación básica entre la probabilidad y la estadística inferencial.



Figura 1.2: Relación básica entre la probabilidad y la estadística inferencial.

Ahora bien, en términos generales, ¿cuál campo es más importante, el de la probabilidad o el de la estadística? Ambos son muy importantes y evidentemente se complementan. La única certeza respecto de la didáctica de ambas disciplinas radica en el hecho de que, si la estadística se debe enseñar con un nivel mayor al de un simple “libro de cocina”, entonces hay que comenzar por enseñar la disciplina de la probabilidad. Esta regla se basa en el hecho de que un analista no podrá aprender nada sobre una población a partir de una muestra hasta que aprenda los rudimentos de incertidumbre en esa muestra. Considere el ejemplo 1.1; en el que la pregunta se centra en si la población, definida por el proceso, tiene o no más de 5% de elementos defectuosos. En otras palabras, la suposición es que 5 de cada 100 artículos, **en promedio, salen defectuosos.** Ahora bien, la muestra contiene 100 artículos y 10 están defectuosos. ¿Esto apoya o refuta la supo-

sición? Aparentemente la refuta porque 10 artículos de cada 100 parecen ser “un trozo grande”. ¿Pero cómo podríamos saber esto sin tener nociones de probabilidad? La única manera en que podremos aprender las condiciones en las cuales el proceso es aceptable (5% de defectuosos) es estudiando el material de los siguientes capítulos. La probabilidad de obtener 10 o más artículos defectuosos en una muestra de 100 es de 0.0282.

Dimos dos ejemplos en donde los elementos de probabilidad ofrecen un resumen que el científico o el ingeniero pueden usar como evidencia para basar una decisión. El puente entre los datos y la conclusión está, por supuesto, basado en los fundamentos de la inferencia estadística, la teoría de la distribución y las distribuciones de muestreos que se examinarán en capítulos posteriores.

1.2 Procedimientos de muestreo; recolección de los datos

En la sección 1.1 estudiamos muy brevemente el concepto de muestreo y el proceso de muestreo. Aunque el muestreo parece ser un concepto simple, la complejidad de las preguntas que se deben contestar acerca de la población, o las poblaciones, en ocasiones requiere que el proceso de muestreo sea muy complejo. El concepto de muestreo se examinará de manera técnica en el capítulo 8, pero aquí nos esforzaremos por dar algunas nociones de sentido común sobre el muestreo. Ésta es una transición natural hacia el análisis del concepto de variabilidad.

Muestreo aleatorio simple

La importancia del muestreo adecuado gira en torno al grado de confianza con que el analista es capaz de responder las preguntas que se plantean. Supongamos que sólo hay una población en el problema. Recuerde que en el ejemplo 1.2 había dos poblaciones implicadas. El **muestreo aleatorio simple** significa que cierta muestra dada de un *tamaño muestral* específico tiene la misma probabilidad de ser seleccionada que cualquiera otra muestra del mismo tamaño. El término **tamaño muestral** simplemente indica el número de elementos en la muestra. Evidentemente, en muchos casos se puede utilizar una tabla de números aleatorios para seleccionar la muestra. La ventaja del muestreo aleatorio simple radica en que ayuda a eliminar el problema de tener una muestra que refleje una población diferente (quizá más restringida) de aquella sobre la cual se necesitan realizar las inferencias. Por ejemplo, se elige una muestra para contestar diferentes preguntas respecto de las preferencias políticas en cierta entidad de Estados Unidos. La muestra implica la elección de, digamos, 1 000 familias y una encuesta a aplicar. Ahora bien, suponga que no se utiliza el muestreo aleatorio, sino que todas o casi todas las 1 000 familias se eligen de una zona urbana. Se considera que las preferencias políticas en las áreas rurales difieren de las de las áreas urbanas. En otras palabras, la muestra obtenida en realidad confinó a la población y, por lo tanto, las inferencias también se tendrán que restringir a la “población confinada”, y en este caso el confinamiento podría resultar indeseable. Si, de hecho, se necesitara hacer las inferencias respecto de la entidad como un todo, a menudo se diría que la muestra con un tamaño de 1 000 familias aquí descrita es una **muestra sesgada**.

Como antes sugerimos, el muestreo aleatorio simple no siempre es adecuado. El enfoque alternativo que se utilice dependerá de la complejidad del problema. Con frecuencia, por ejemplo, las unidades muestrales no son homogéneas y se dividen naturalmente en grupos que no se traslapan y que son homogéneos. Tales grupos se llaman *estratos*, y

un procedimiento llamado *muestreo aleatorio estratificado* implica la selección al azar de una muestra *dentro* de cada estrato. El propósito de esto es asegurarse de que ninguno de los estratos esté sobrerrepresentado ni subrepresentado. Por ejemplo, suponga que se aplica una encuesta a una muestra para reunir opiniones preliminares respecto de un referéndum que se piensa realizar en determinada ciudad. La ciudad está subdividida en varios grupos étnicos que representan estratos naturales y, para no excluir ni sobrerrepresentar a algún grupo de cada uno de ellos, se eligen muestras aleatorias separadas de cada grupo.

Diseño experimental

El concepto de aleatoriedad o asignación aleatoria desempeña un papel muy importante en el área del **diseño experimental**, que se presentó brevemente en la sección 1.1 y es un fundamento muy importante en casi cualquier área de la ingeniería y de la ciencia experimental. Estudiaremos este tema con detenimiento en los capítulos 13 a 15. Sin embargo, es conveniente introducirlo aquí brevemente en el contexto del muestreo aleatorio. Un conjunto de los llamados **tratamientos** o **combinaciones de tratamientos** se vuelven las poblaciones que se van a estudiar o a comparar en algún sentido. Un ejemplo es el tratamiento “con nitrógeno” *versus* “sin nitrógeno” del ejemplo 1.2. Otro ejemplo sencillo sería “placebo” *versus* “medicamento activo” o, en un estudio sobre la fatiga por corrosión, tendríamos combinaciones de tratamientos que impliquen especímenes con recubrimiento o sin recubrimiento, así como condiciones de alta o de baja humedad, a las cuales se somete el espécimen. De hecho, habrían cuatro combinaciones de factores o de tratamientos (es decir, 4 poblaciones), y se podrían formular y responder muchas preguntas científicas usando los métodos estadísticos e inferenciales. Considere primero la situación del ejemplo 1.2. En el experimento hay 20 plantones enfermos implicados. A partir de los datos es fácil observar que los plantones son diferentes entre sí. Dentro del grupo tratado con nitrógeno (o del grupo que no se trató con nitrógeno) hay **variabilidad** considerable en el peso de los tallos, la cual se debe a lo que por lo general se denomina **unidad experimental**. Éste es un concepto tan importante en la estadística inferencial que no es posible describirlo totalmente en este capítulo. La naturaleza de la variabilidad es muy importante. Si es demasiado grande, debido a que resulta de una condición de excesiva falta de homogeneidad en las unidades experimentales, la variabilidad “eliminará” cualquier diferencia detectable entre ambas poblaciones. Recuerde que en este caso eso no ocurrió.

La gráfica de puntos de la figura 1.1 y el valor-*P* indican una clara distinción entre esas dos condiciones. ¿Qué papel desempeñan tales unidades experimentales en el proceso mismo de recolección de los datos? El enfoque por sentido común y, de hecho, estándar, es asignar los 20 plantones o unidades experimentales **aleatoriamente a las dos condiciones o tratamientos**. En el estudio del medicamento podríamos decidir utilizar un total de 200 pacientes disponibles, quienes serán claramente distinguibles en algún sentido. Ellos son las unidades experimentales. No obstante, tal vez todos tengan una condición crónica que podría ser tratada con el fármaco. Así, en el denominado **diseño completamente aleatorio**, se asignan al azar 100 pacientes al placebo y 100 al medicamento activo. De nuevo, son estas unidades experimentales en el grupo o tratamiento las que producen la variabilidad en el resultado de los datos (es decir, la variabilidad en el resultado medido), digamos, de la presión sanguínea o cualquier valor de la eficacia de un medicamento que sea importante. En el estudio de la fatiga por corrosión las unidades experimentales son los especímenes que se someten a la corrosión.

¿Por qué las unidades experimentales se asignan aleatoriamente?

¿Cuál es el posible efecto negativo de no asignar aleatoriamente las unidades experimentales a los tratamientos o a las combinaciones de tratamientos? Esto se observa más claramente en el caso del estudio del medicamento. Entre las características de los pacientes que producen variabilidad en los resultados están la edad, el género y el peso. Tan sólo suponga que por casualidad el grupo del placebo contiene una muestra de personas que son predominantemente más obesas que las del grupo del tratamiento. Quizá los individuos más obesos muestren una tendencia a tener una presión sanguínea más elevada, lo cual evidentemente sesgará el resultado y, por lo tanto, cualquier resultado que se obtenga al aplicar la inferencia estadística podría tener poco que ver con el efecto del medicamento, pero mucho con las diferencias en el peso de ambas muestras de pacientes.

Deberíamos enfatizar la importancia del término **variabilidad**. La variabilidad excesiva entre las unidades experimentales “disfraza” los hallazgos científicos. En secciones posteriores intentaremos clasificar y cuantificar las medidas de variabilidad. En las siguientes secciones presentaremos y analizaremos cantidades específicas que se calculan en las muestras; las cantidades proporcionan una idea de la naturaleza de la muestra respecto de la ubicación del centro de los datos y la variabilidad de los mismos. Un análisis de varias de tales medidas de un solo número permite ofrecer un preámbulo de que la información estadística será un componente importante de los métodos estadísticos que se utilizarán en capítulos posteriores. Estas medidas, que ayudan a clasificar la naturaleza del conjunto de datos, caen en la categoría de **estadísticas descriptivas**. Este material es una introducción a una presentación breve de los métodos pictóricos y gráficos que van incluso más allá en la caracterización del conjunto de datos. El lector debería entender que los métodos estadísticos que se presentan aquí se utilizarán a lo largo de todo el texto. Para ofrecer una imagen más clara de lo que implican los estudios de diseño experimental se presenta el ejemplo 1.3.

Ejemplo 1.3: Se realizó un estudio sobre la corrosión con la finalidad de determinar si al recubrir una aleación de aluminio con una sustancia retardadora de la corrosión, el metal se corroe menos. El recubrimiento es un protector que los anunciantes afirman que minimiza el daño por fatiga en esta clase de material. La influencia de la humedad sobre la magnitud de la corrosión también es de interés. Una medición de la corrosión puede expresarse en millares de ciclos hasta la ruptura del metal. Se utilizaron dos niveles de recubrimiento: sin recubrimiento y con recubrimiento químico contra la corrosión. También se consideraron dos niveles de humedad relativa, de 20% y 80%, respectivamente.

El experimento implica las cuatro combinaciones de tratamientos que se listan en la siguiente tabla. Se usan ocho unidades experimentales, que son especímenes de aluminio preparados, dos de los cuales se asignan aleatoriamente a cada una de las cuatro combinaciones de tratamiento. Los datos se presentan en la tabla 1.2.

Los datos de la corrosión son promedios de los dos especímenes. En la figura 1.3 se presenta una gráfica con los promedios. Un valor relativamente grande de ciclos hasta la ruptura representa una cantidad pequeña de corrosión. Como se podría esperar, al parecer un incremento en la humedad hace que empeore la corrosión. El uso del procedimiento de recubrimiento químico contra la corrosión parece reducir la corrosión. ■

En este ejemplo de diseño experimental el ingeniero eligió sistemáticamente las cuatro combinaciones de tratamiento. Para vincular esta situación con los conceptos con los que el lector se ha familiarizado hasta aquí, deberíamos suponer que las condiciones

Tabla 1.2: Datos para el ejemplo 1.3

Recubrimiento	Humedad	Promedio de corrosión
		en miles de ciclos hasta la ruptura
Sin recubrimiento	20%	975
	80%	350
Con recubrimiento químico contra la corrosión	20%	1750
	80%	1550

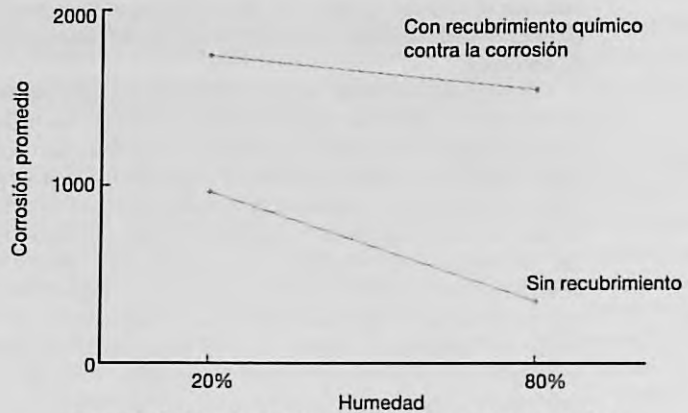


Figura 1.3: Resultados de corrosión para el ejemplo 1.3.

que representan las cuatro combinaciones de tratamientos son cuatro poblaciones separadas y que los dos valores de corrosión observados en cada una de las poblaciones constituyen importantes piezas de información. La importancia del promedio al captar y resumir ciertas características en la población se destacará en la sección 1.3. Aunque a partir de la figura podríamos sacar conclusiones acerca del papel que desempeña la humedad y del efecto de recubrir los especímenes, no podemos evaluar con exactitud los resultados de un punto de vista analítico sin tomar en cuenta la *variabilidad alrededor del promedio*. De nuevo, como señalamos con anterioridad, si los dos valores de corrosión en cada una de las combinaciones de tratamientos son muy cercanos, la imagen de la figura 1.3 podría ser una descripción precisa. Pero si cada valor de la corrosión en la figura es un promedio de dos valores que están ampliamente dispersos, entonces esta variabilidad podría, de hecho, en verdad “eliminar” cualquier información que parezca difundirse cuando tan sólo se observan los promedios. Los siguientes ejemplos ilustran estos conceptos:

1. La asignación aleatoria a las combinaciones de tratamientos (recubrimiento/humedad) de las unidades experimentales (especímenes).
2. El uso de promedios muestrales (valores de corrosión promedio) para resumir la información muestral.
3. La necesidad de considerar las medidas de variabilidad en el análisis de cualquier muestra o conjunto de muestras.

Este ejemplo sugiere la necesidad de estudiar el tema que se expone en las secciones 1.3 y 1.4, es decir, el de las estadísticas descriptivas que indican las medidas de la ubicación del centro en un conjunto de datos, y aquellas con las que se mide la variabilidad.

1.3 Medidas de localización: la media y la mediana de una muestra

Las medidas de localización están diseñadas para brindar al analista algunos valores cuantitativos de la ubicación central o de otro tipo de los datos en una muestra. En el ejemplo 1.2 parece que el centro de la muestra con nitrógeno claramente excede al de la muestra sin nitrógeno. Una medida obvia y muy útil es la **media de la muestra**. La media es simplemente un promedio numérico.

Definición 1.1: Suponga que las observaciones en una muestra son x_1, x_2, \dots, x_n . La **media de la muestra**, que se denota con \bar{x} , es

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}.$$

Hay otras medidas de tendencia central que se explican con detalle en capítulos posteriores. Una medida importante es la **mediana de la muestra**. El propósito de la mediana de la muestra es reflejar la tendencia central de la muestra de manera que no sea influida por los valores extremos.

Definición 1.2: Dado que las observaciones en una muestra son x_1, x_2, \dots, x_n , acomodadas en **orden de magnitud creciente**, la mediana de la muestra es

$$\tilde{x} = \begin{cases} x_{(n+1)/2}, & \text{si } n \text{ es impar,} \\ \frac{1}{2}(x_{n/2} + x_{n/2+1}), & \text{si } n \text{ es par.} \end{cases}$$

Por ejemplo, suponga que el conjunto de datos es el siguiente: 1.7, 2.2, 3.9, 3.11 y 14.7. La media y la mediana de la muestra son, respectivamente,

$$\bar{x} = 5.12, \quad \tilde{x} = 3.9.$$

Es evidente que la media es influida de manera considerable por la presencia de la observación extrema, 14.7; en tanto que el lugar de la mediana hace énfasis en el verdadero "centro" del conjunto de datos. En el caso del conjunto de datos de dos muestras del ejemplo 1.2, las dos medidas de tendencia central para las muestras individuales son

$$\begin{aligned} \bar{x} \text{ (sin nitrógeno)} &= 0.399 \text{ gramos,} \\ \bar{x} \text{ (sin nitrógeno)} &= \frac{0.38 + 0.42}{2} = 0.400 \text{ gramos,} \\ \bar{x} \text{ (con nitrógeno)} &= 0.565 \text{ gramos,} \\ \bar{x} \text{ (con nitrógeno)} &= \frac{0.49 + 0.52}{2} = 0.505 \text{ gramos.} \end{aligned}$$

Es evidente que hay una diferencia conceptual entre la media y la mediana. Para el lector con ciertas nociones de ingeniería quizá sea de interés que la media de la muestra

es el **centroide de los datos** en una muestra. En cierto sentido es el punto en el cual se puede colocar un fulcro (apoyo) para equilibrar un sistema de “pesos”, que son las ubicaciones de los datos individuales. Esto se muestra en la figura 1.4 respecto de la muestra “con nitrógeno”.

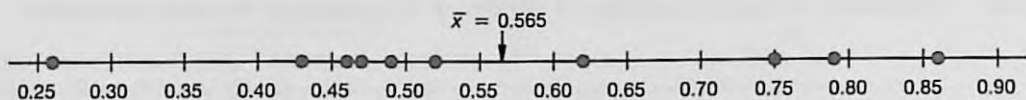


Figura 1.4: Media de la muestra como centroide del peso del tallo con nitrógeno.

En capítulos posteriores la base para el cálculo de \bar{x} es un **estimado de la media de la población**. Como antes señalamos, el propósito de la inferencia estadística es obtener conclusiones acerca de las características o **parámetros** y la **estimación** es una característica muy importante de la inferencia estadística.

La mediana y la media pueden ser muy diferentes entre sí. Observe, sin embargo, que en el caso de los datos del peso de los tallos el valor de la media de la muestra para “sin nitrógeno” es bastante similar al valor de la mediana.

Otras medidas de localización

Hay muchos otros métodos para calcular la ubicación del centro de los datos en la muestra. No los trataremos en este momento. Por lo general las alternativas para la media de la muestra se diseñan con el fin de generar valores que representen relación entre la media y la mediana. Rara vez utilizamos alguna de tales medidas. Sin embargo, es aleccionador estudiar una clase de estimadores conocida como **media recortada**, la cual se calcula “quitando” cierto porcentaje de los valores mayores y menores del conjunto. Por ejemplo, la media recortada al 10% se encuentra eliminando tanto el 10% de los valores mayores como el 10% de los menores, y calculando el promedio de los valores restantes. En el caso de los datos del peso de los tallos, eliminaríamos el valor más alto y el más bajo, ya que el tamaño de la muestra es 10 en cada caso. De manera que para el grupo sin nitrógeno la media recortada al 10% está dada por

$$\bar{x}_{\text{rec}(10)} = \frac{0.32 + 0.37 + 0.47 + 0.43 + 0.36 + 0.42 + 0.38 + 0.43}{8} = 0.39750,$$

y para la media recortada al 10% del grupo con nitrógeno tenemos

$$\bar{x}_{\text{rec}(10)} = \frac{0.43 + 0.47 + 0.49 + 0.52 + 0.75 + 0.79 + 0.62 + 0.46}{8} = 0.56625.$$

Observe que en este caso, como se esperaba, las medias recortadas están cerca tanto de la media como de la mediana para las muestras individuales. Desde luego, el enfoque de la media recortada es menos sensible a los valores extremos que la media de la muestra, pero no tan insensible como la mediana. Además, el método de la media recortada utiliza más información que la mediana de la muestra. Advierta que la mediana de la muestra es, de hecho, un caso especial de la media recortada, en el cual se eliminan todos los datos de la muestra y queda sólo el central o dos observaciones.

Ejercicios

1.1 Se registran las siguientes mediciones para el tiempo de secado (en horas) de cierta marca de pintura esmaltada.

3.4	2.5	4.8	2.9	3.6
2.8	3.3	5.6	3.7	2.8
4.4	4.0	5.2	3.0	4.8

Suponga que las mediciones constituyen una muestra aleatoria simple.

- ¿Cuál es el tamaño de la muestra anterior?
- Calcule la media de la muestra para estos datos.
- Calcule la mediana de la muestra.
- Grafique los datos utilizando una gráfica de puntos.
- Calcule la media recortada al 20% para el conjunto de datos anterior.
- ¿La media muestral para estos datos es más o menos descriptiva como centro de localización, que la media recortada?

1.2 Según la revista *Chemical Engineering*, una propiedad importante de una fibra es su absorción del agua. Se toma una muestra aleatoria de 20 pedazos de fibra de algodón y se mide la absorción de cada uno. Los valores de absorción son los siguientes:

18.71	21.41	20.72	21.81	19.29	22.43	20.17
23.71	19.44	20.50	18.92	20.33	23.00	22.85
19.25	21.77	22.11	19.77	18.04	21.12	

- Calcule la media y la mediana muestrales para los valores de la muestra anterior.
- Calcule la media recortada al 10%.
- Elabore una gráfica de puntos con los datos de la absorción.
- Si se utilizan sólo los valores de la media, la mediana y la media recortada, ¿hay evidencia de valores extremos en los datos?

1.3 Se utiliza cierto polímero para los sistemas de evacuación de los aviones. Es importante que el polímero sea resistente al proceso de envejecimiento. Se utilizaron veinte especímenes del polímero en un experimento. Diez se asignaron aleatoriamente para exponerse a un proceso de envejecimiento acelerado del lote, el cual implica la exposición a altas temperaturas durante 10 días. Se hicieron las mediciones de resistencia a la tensión de los especímenes y se registraron los siguientes datos sobre resistencia a la tensión en psi.

Sin envejecimiento:	227	222	218	217	225
	218	216	229	228	221
Con envejecimiento:	219	214	215	211	209
	218	203	204	201	205

- Elabore la gráfica de puntos de los datos.
- ¿En la gráfica que obtuvo parece que el proceso de envejecimiento tuvo un efecto en la resistencia

a la tensión de este polímero? Explique su respuesta.

- Calcule la resistencia a la tensión de la media de la muestra en las dos muestras.
- Calcule la mediana de ambas. Analice la similitud o falta de similitud entre la media y la mediana de cada grupo.

1.4 En un estudio realizado por el Departamento de Ingeniería Mecánica del Tecnológico de Virginia se compararon las varillas de acero que abastecen dos compañías diferentes. Se fabricaron diez resortes de muestra con las varillas de metal proporcionadas por cada una de las compañías y se registraron sus medidas de flexibilidad. Los datos son los siguientes:

Compañía A:	9.3	8.8	6.8	8.7	8.5
	6.7	8.0	6.5	9.2	7.0
Compañía B:	11.0	9.8	9.9	10.2	10.1
	9.7	11.0	11.1	10.2	9.6

- Calcule la media y la mediana de la muestra para los datos de ambas compañías.
- Grafique los datos para las dos compañías en la misma línea y explique su conclusión respecto de cualquier aparente diferencia entre las dos compañías.

1.5 Veinte hombres adultos de entre 30 y 40 años de edad participaron en un estudio para evaluar el efecto de cierto régimen de salud, que incluye dieta y ejercicio, en el colesterol sanguíneo. Se eligieron aleatoriamente diez para el grupo de control y los otros diez se asignaron para participar en el régimen como el grupo de tratamiento durante un periodo de seis meses. Los siguientes datos muestran la reducción en el colesterol que experimentaron en ese periodo los 20 sujetos:

Grupo de control:	7	3	-4	14	2
	5	22	-7	9	5
Grupo de tratamiento:	-6	5	9	4	4
	12	37	5	3	3

- Elabore una gráfica de puntos con los datos de ambos grupos en la misma gráfica.
- Calcule la media, la mediana y la media recortada al 10% para ambos grupos.
- Explique por qué la diferencia en las medias sugiere una conclusión acerca del efecto del régimen, en tanto que la diferencia en las medianas o las medias recortadas sugiere una conclusión diferente.

1.6 La resistencia a la tensión del caucho de silicio se considera una función de la temperatura de vulcanizado. Se llevó a cabo un estudio en el que se prepararon muestras de 12 especímenes del caucho utilizando temperaturas de vulcanizado de 20°C y 45°C. Los siguientes

Rango y desviación estándar de la muestra

Así como hay muchas medidas de tendencia central o de localización, hay muchas medidas de dispersión o variabilidad. Quizá la más simple sea el **rango de la muestra** $X_{\max} - X_{\min}$. El rango puede ser muy útil y se examina con amplitud en el capítulo 17 sobre *control estadístico de calidad*. La medida muestral de dispersión que se utiliza más a menudo es la **desviación estándar de la muestra**. Nuevamente denotemos con x_1, x_2, \dots, x_n los valores de la muestra.

Definición 1.3: La **varianza de la muestra**, denotada con s^2 , está dada por

$$s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}.$$

La **desviación estándar de la muestra**, denotada con s , es la raíz cuadrada positiva de s^2 , es decir,

$$s = \sqrt{s^2}.$$

Para el lector debería quedar claro que la desviación estándar de la muestra es, de hecho, una medida de variabilidad. Una variabilidad grande en un conjunto de datos produce valores relativamente grandes de $(x - \bar{x})^2$ y, por consiguiente, una varianza muestral grande. La cantidad $n - 1$ a menudo se denomina **grados de libertad asociados con la varianza** estimada. En este ejemplo sencillo los grados de libertad representan el número de piezas de información independientes disponibles para calcular la variabilidad. Por ejemplo, suponga que deseamos calcular la varianza de la muestra γ la desviación estándar del conjunto de datos (5, 17, 6, 4). El promedio de la muestra es $\bar{x} = 8$. El cálculo de la varianza implica:

$$(5 - 8)^2 + (17 - 8)^2 + (6 - 8)^2 + (4 - 8)^2 = (-3)^2 + 9^2 + (-2)^2 + (-4)^2.$$

Las cantidades dentro de los paréntesis suman cero. En general, $\sum_{i=1}^n (x_i - \bar{x}) = 0$ (véase el ejercicio 1.16 de la página 31). Entonces, el cálculo de la varianza de una muestra no implica n **desviaciones cuadradas independientes** de la media \bar{x} . De hecho, como el último valor de $x - \bar{x}$ es determinado por los primeros $n - 1$ valores, decimos que éstas son $n - 1$ “piezas de información” que produce s^2 . Por consiguiente, hay $n - 1$ grados de libertad en vez de n grados de libertad para calcular la varianza de una muestra.

Ejemplo 1.4: En un ejemplo que se estudia ampliamente en el capítulo 10, un ingeniero se interesa en probar el “sesgo” en un medidor de pH. Los datos se recaban con el medidor mediante la medición del pH de una sustancia neutra (pH = 7.0). Se toma una muestra de tamaño 10 y se obtienen los siguientes resultados:

7.07 7.00 7.10 6.97 7.00 7.03 7.01 7.01 6.98 7.08.

La media de la muestra \bar{x} está dada por

$$\bar{x} = \frac{7.07 + 7.00 + 7.10 + \dots + 7.08}{10} = 7.0250.$$

La varianza de la muestra s^2 está dada por

$$s^2 = \frac{1}{9}[(7.07 - 7.025)^2 + (7.00 - 7.025)^2 + (7.10 - 7.025)^2 + \dots + (7.08 - 7.025)^2] = 0.001939.$$

Como resultado, la desviación estándar de la muestra está dada por

$$s = \sqrt{0.001939} = 0.044.$$

Así que la desviación estándar de la muestra es 0.0440 con $n - 1 = 9$ grados de libertad. ■

Unidades para la desviación estándar y la varianza

A partir de la definición 1.3 debería ser evidente que la varianza es una medida de la desviación cuadrática promedio de la media \bar{x} . Empleamos el término *desviación cuadrática promedio* aun cuando la definición utilice una división entre $n - 1$ grados de libertad en vez de n . Desde luego, si n es grande, la diferencia en el denominador es inconsecuente. Por lo tanto, la varianza de la muestra tiene unidades que son el cuadrado de las unidades en los datos observados; aunque la desviación estándar de la muestra se encuentra en unidades lineales. Considere los datos del ejemplo 1.2. Los pesos del tallo se miden en gramos. Como resultado, las desviaciones estándar de la muestra están en gramos y las varianzas se miden en gramos². De hecho, las desviaciones estándar individuales son 0.0728 gramos para el caso sin nitrógeno y 0.1867 gramos para el grupo con nitrógeno. Observe que la desviación estándar en verdad indica una variabilidad mucho más grande en la muestra con nitrógeno. Esta condición se destaca en la figura 1.1.

¿Cuál es la medida de variabilidad más importante?

Como indicamos antes, el rango de la muestra tiene aplicaciones en el área del control estadístico de la calidad. Quizás el lector considere que es redundante utilizar la varianza de la muestra y la desviación estándar de la muestra. Ambas medidas reflejan el mismo concepto en la variabilidad de la medición, pero la desviación estándar de la muestra mide la variabilidad en unidades lineales; en tanto que la varianza muestral se mide en unidades cuadradas. Ambas desempeñan papeles importantes en el uso de los métodos estadísticos. Mucho de lo que se logra en el contexto de la inferencia estadística implica la obtención de conclusiones acerca de las características de poblaciones. Entre tales características son constantes los denominados **parámetros de la población**. Dos parámetros importantes son la **media de la población** y la **varianza de la población**. La varianza de la muestra desempeña un papel explícito en los métodos estadísticos que se utilizan para obtener inferencias sobre la varianza de la población. La desviación estándar de la muestra desempeña un papel importante, junto con la media de la muestra, en las inferencias que se realizan acerca de la media de la población. En general, la varianza se considera más en la teoría inferencial, mientras que la desviación estándar se utiliza más en aplicaciones.

Ejercicios

1.7 Considere los datos del tiempo de secado del ejercicio 1.1 de la página 13. Calcule la varianza de la muestra y la desviación estándar de la muestra.

1.8 Calcule la varianza de la muestra y la desviación estándar para los datos de absorción del agua del ejercicio 1.2 de la página 13.

1.9 El ejercicio 1.3 de la página 13 presentó datos de resistencia a la tensión de dos muestras, una en la que los especímenes se expusieron a un proceso de envejecimiento y otra en la que no se efectuó tal proceso en los especímenes.

- Calcule la varianza de la muestra, así como su desviación estándar, en cuanto a la resistencia a la tensión en ambas muestras.
- ¿Parece haber alguna evidencia de que el envejecimiento afecta la variabilidad en la resistencia a la

tensión? (Véase también la gráfica para el ejercicio 1.3 de la página 13).

1.10 Para los datos del ejercicio 1.4 de la página 13 calcule tanto la media como la varianza de la "flexibilidad" para las compañías A y B. ¿Parece que hay una diferencia de flexibilidad entre la compañía A y la compañía B?

1.11 Considere los datos del ejercicio 1.5 de la página 13. Calcule la varianza de la muestra y la desviación estándar de la muestra para ambos grupos: el de tratamiento y el de control.

1.12 Para el ejercicio 1.6 de la página 13 calcule la desviación estándar muestral de la resistencia a la tensión para las muestras, de forma separada para ambas temperaturas. ¿Parece que un incremento en la temperatura influye en la variabilidad de la resistencia a la tensión? Explique su respuesta.

1.5 Datos discretos y continuos

La inferencia estadística a través del análisis de estudios observacionales o de diseños experimentales se utiliza en muchas áreas científicas. Los datos reunidos pueden ser **discretos** o **continuos**, según el área de aplicación. Por ejemplo, un ingeniero químico podría estar interesado en un experimento que lo lleve a condiciones en que se maximice la producción. Aquí, por supuesto, la producción se expresaría en porcentaje, o gramos/libra, medida en un continuo. Por otro lado, un toxicólogo que realice un experimento de combinación de fármacos quizás encuentre datos que son binarios por naturaleza (es decir, el paciente responde o no lo hace).

En la teoría de la probabilidad se hacen distinciones importantes entre datos discretos y continuos que nos permiten hacer inferencias estadísticas. Con frecuencia las aplicaciones de la inferencia estadística se encuentran cuando se trabaja con *datos por conteo*. Por ejemplo, un ingeniero podría estar interesado en estudiar el número de partículas radiactivas que pasan a través de un contador en, digamos, 1 milisegundo. Al personal responsable de la eficiencia de una instalación portuaria podría interesarle conocer las características del número de buques petroleros que llegan diariamente a cierta ciudad portuaria. En el capítulo 5 se examinarán varios escenarios diferentes que conducen a distintas formas de manejo de los datos para situaciones con datos por conteo.

Incluso en esta fase inicial del texto se debería poner especial atención a algunos detalles que se asocian con datos binarios. Son muchas las aplicaciones que requieren el análisis estadístico de datos binarios. Con frecuencia la medición que se utiliza en el análisis es la *proporción muestral*. En efecto, la situación binaria implica dos categorías. Si en los datos hay n unidades y x se define como el número que cae en la categoría 1, entonces $n - x$ cae en la categoría 2. Así, x/n es la proporción muestral en la categoría 1 y $1 - x/n$ es la proporción muestral en la categoría 2. En la aplicación biomédica, por ejemplo, 50 pacientes representarían las unidades de la muestra y si, después de que se les suministra el medicamento, 20 de los 50 experimentarían mejoría en sus malestares estomacales (que son comunes en los 50), entonces $\frac{20}{50} = 0.4$ sería la proporción muestral

para la cual el medicamento tuvo éxito, y $1 - 0.4 = 0.6$ sería la proporción muestral para la cual el fármaco no tuvo éxito. En realidad, la medición numérica fundamental para datos binarios por lo general se denota con 0 o con 1. Éste es el caso de nuestro ejemplo médico, en el que un resultado exitoso se denota con un 1 y uno no exitoso con un 0. Entonces, la proporción muestral es en realidad una media muestral de unos y ceros. Para la categoría de éxitos,

$$\frac{x_1 + x_2 + \cdots + x_{50}}{50} = \frac{1 + 1 + 0 + \cdots + 0 + 1}{50} = \frac{20}{50} = 0.4.$$

¿Qué clases de problemas se resuelven en situaciones con datos binarios?

Los tipos de problemas que enfrentan científicos e ingenieros que usan datos binarios no son muy difíciles, a diferencia de aquellos en los que las mediciones de interés son las continuas. Sin embargo, se utilizan técnicas diferentes debido a que las propiedades estadísticas de las proporciones muestrales son bastante diferentes de las medias muestrales que resultan de los promedios tomados de poblaciones continuas. Considere los datos del ejemplo en el ejercicio 1.6 de la página 13. El problema estadístico subyacente en este caso se enfoca en si una intervención, digamos un incremento en la temperatura de vulcanizado, alterará la resistencia a la tensión de la media de la población que se asocia con el proceso del caucho de silicio. Por otro lado, en el área de control de calidad, suponga que un fabricante de neumáticos para automóvil informa que en un embarque con 5000 neumáticos, seleccionados aleatoriamente del proceso, hay 100 defectuosos. Aquí la proporción muestral es $\frac{100}{5000} = 0.02$. Luego de realizar un cambio en el proceso diseñado para reducir los neumáticos defectuosos, se toma una segunda muestra de 5000 y se encuentran 90 defectuosos. La proporción muestral se redujo a $\frac{90}{5000} = 0.018$. Entonces, surge una pregunta: “¿La disminución en la proporción muestral de 0.02 a 0.018 es suficiente para sugerir una mejoría real en la proporción de la población?” En ambos casos se requiere el uso de las propiedades estadísticas de los promedios de la muestra: uno de las muestras de poblaciones continuas y el otro de las muestras de poblaciones discretas (binarias). En ambos casos la media de la muestra es un **estimado** de un parámetro de la población, una media de la población en el primer caso (es decir, la media de la resistencia a la tensión) y una proporción de la población (o sea, la proporción de neumáticos defectuosos en la población) en el segundo caso. Así que aquí tenemos estimados de la muestra que se utilizan para obtener conclusiones científicas respecto de los parámetros de la población. Como indicamos en la sección 1.3, éste es el tema general en muchos problemas prácticos en los que se usa la inferencia estadística.

1.6 Modelado estadístico, inspección científica y diagnósticos gráficos

A menudo, el resultado final de un análisis estadístico es la estimación de los parámetros de un **modelo postulado**. Éste es un proceso natural para los científicos y los ingenieros, ya que con frecuencia usan modelos. Un modelo estadístico no es determinista, es más bien un modelo que conlleva algunos aspectos probabilísticos. A menudo una forma de modelo es la base de las **suposiciones** que hace el analista. En el ejemplo 1.2 el científico podría desear determinar, a través de la información de la muestra, algún nivel de distinción entre las poblaciones tratadas con nitrógeno y las poblaciones no tratadas. El análisis podría requerir cierto modelo para los datos; por ejemplo, que las dos muestras

proviengan de **distribuciones normales o gaussianas**. Véase el capítulo 6 para el estudio de la distribución normal.

Es evidente que quienes utilizan métodos estadísticos no pueden generar la información o los datos experimentales suficientes para describir a la totalidad de la población. Pero es frecuente que se utilicen los conjuntos de datos para aprender sobre ciertas propiedades de la población. Los científicos y los ingenieros están acostumbrados a manejar conjuntos de datos. Debería ser obvia la importancia de describir o *resumir* la naturaleza de los conjuntos de datos. Con frecuencia el resumen gráfico de un conjunto de datos puede proporcionar información sobre el sistema del que se obtuvieron los datos. Por ejemplo, en las secciones 1.1 y 1.3 mostramos gráficas de puntos.

En esta sección se estudia con detalle el papel del muestreo y de la graficación de los datos para mejorar la **inferencia estadística**. Nos limitamos a presentar algunas gráficas sencillas, pero a menudo efectivas, que complementan el estudio de poblaciones estadísticas.

Diagrama de dispersión

A veces el modelo postulado puede tener una forma algo más compleja. Por ejemplo, considere a un fabricante de textiles que diseña un experimento en donde se producen especímenes de tela que contienen diferentes porcentajes de algodón. Considere los datos de la tabla 1.3.

Tabla 1.3: Resistencia a la tensión

Porcentaje del algodón	Resistencia a la tensión
15	7, 7, 9, 8, 10
20	19, 20, 21, 20, 22
25	21, 21, 17, 19, 20
30	8, 7, 8, 9, 10

Se fabrican cinco especímenes de tela para cada uno de los cuatro porcentajes de algodón. En este caso tanto el modelo para el experimento como el tipo de análisis que se utiliza deberían tomar en cuenta el objetivo del experimento y los insumos importantes del científico textil. Algunas gráficas sencillas podrían mostrar la clara distinción entre las muestras. Véase la figura 1.5; las medias y la variabilidad muestrales se describen bien en el diagrama de dispersión. El objetivo de este experimento podría ser simplemente determinar cuáles porcentajes de algodón son verdaderamente distintos de los otros. En otras palabras, como en el caso de los datos con nitrógeno y sin nitrógeno, ¿para cuáles porcentajes de algodón existen diferencias claras entre las poblaciones o, de forma más específica, entre las medias de las poblaciones? En este caso quizás un modelo razonable es que cada muestra proviene de una distribución normal. Aquí el objetivo es muy semejante al de los datos con nitrógeno y sin nitrógeno, excepto que se incluyen más muestras. El formalismo del análisis implica nociones de prueba de hipótesis, los cuales se examinarán en el capítulo 10. A propósito, quizás este formalismo no sea necesario a la luz del diagrama de diagnóstico. Pero, ¿describe éste el objetivo real del experimento y, por consiguiente, el enfoque adecuado para el análisis de datos? Es probable que el científico anticipe la existencia de una *resistencia a la tensión máxima de la media de la población* en el rango de concentración de algodón en el experimento. Aquí el análisis de los datos debería girar en torno a un tipo diferente de modelo, es decir, uno

que postule un tipo de estructura que relacione la resistencia a la tensión de la media de la población con la concentración de algodón. En otras palabras, un modelo se puede escribir como

$$\mu_{t,c} = \beta_0 + \beta_1 C + \beta_2 C^2,$$

en donde $\mu_{t,c}$ es la resistencia a la tensión de la media de la población, que varía con la cantidad de algodón en el producto C . La implicación de este modelo es que, para un nivel fijo de algodón, hay una población de mediciones de resistencia a la tensión y la media de la población es $\mu_{t,c}$. Este tipo de modelo, que se denomina **modelo de regresión**, se estudiará en los capítulos 11 y 12. La forma funcional la elige el científico. A veces el análisis de datos puede sugerir que se cambie el modelo. Entonces el analista de datos “considera” un modelo que se pueda alterar después de hacer cierto análisis. El uso de un modelo empírico va acompañado por la **teoría de estimación**, donde β_0 , β_1 y β_2 se estiman a partir de los datos. Además, la inferencia estadística se puede, entonces, utilizar para determinar lo adecuado del modelo.

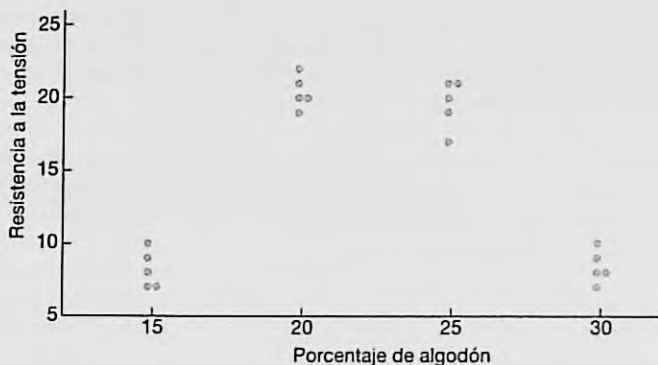


Figura 1.5: Diagrama de dispersión de la resistencia a la tensión y los porcentajes de algodón.

Aquí se hacen evidentes dos puntos de las dos ilustraciones de datos: 1) el tipo de modelo que se emplea para describir los datos a menudo depende del objetivo del experimento, y 2) la estructura del modelo debería aprovechar el insumo científico no estadístico. La selección de un modelo representa una **suposición fundamental** sobre la que se basa la inferencia estadística resultante. A lo largo del libro se hará evidente la importancia que las gráficas pueden llegar a tener. A menudo las gráficas ilustran información que permite que los resultados de la inferencia estadística formal se comuniquen mejor al científico o al ingeniero. A veces las gráficas o el **análisis exploratorio de los datos** pueden enseñar al analista información que no se obtiene del análisis formal. Casi cualquier análisis formal requiere suposiciones que se desarrollan a partir del modelo de datos. Las gráficas pueden resaltar la **violación de suposiciones** que de otra forma no se notarían. A lo largo del libro las gráficas se utilizarán de manera extensa para complementar el análisis formal de los datos. En las siguientes secciones se presentan algunas herramientas gráficas que son útiles para el análisis exploratorio o descriptivo de los datos.

Diagrama de tallo y hojas

Los datos estadísticos obtenidos de poblaciones grandes pueden ser muy útiles para estudiar el comportamiento de la distribución si se presentan en una combinación tabular y gráfica conocida como **diagrama de tallo y hojas**.

Para ejemplificar la elaboración de un diagrama de tallo y hojas considere los datos de la tabla 1.4, que especifican la “vida” de 40 baterías para automóvil similares, registradas al décimo de año más cercano. Las baterías se garantizan por tres años. Comience por dividir cada observación en dos partes: una para el tallo y otra para las hojas, de manera que el tallo represente el dígito entero que antecede al decimal y la hoja corresponda a la parte decimal del número. En otras palabras, para el número 3.7 el dígito 3 se designa al tallo y el 7 a la hoja. Para nuestros datos los cuatro tallos 1, 2, 3 y 4 se listan verticalmente del lado izquierdo de la tabla 1.5, en tanto que las hojas se registran en el lado derecho correspondiente al valor del tallo adecuado. Entonces, la hoja 6 del número 1.6 se registra enfrente del tallo 1; la hoja 5 del número 2.5 enfrente del tallo 2; y así sucesivamente. El número de hojas registrado junto a cada tallo se anota debajo de la columna de frecuencia.

Tabla 1.4: Vida de las baterías para automóvil

2.2	4.1	3.5	4.5	3.2	3.7	3.0	2.6
3.4	1.6	3.1	3.3	3.8	3.1	4.7	3.7
2.5	4.3	3.4	3.6	2.9	3.3	3.9	3.1
3.3	3.1	3.7	4.4	3.2	4.1	1.9	3.4
4.7	3.8	3.2	2.6	3.9	3.0	4.2	3.5

Tabla 1.5: Diagrama de tallo y hojas de la vida de las baterías

Tallo	Hoja	Frecuencia
1	69	2
2	25669	5
3	0011112223334445567778899	25
4	11234577	8

El diagrama de tallo y hojas de la tabla 1.5 contiene sólo cuatro tallos y, en consecuencia, no ofrece una representación adecuada de la distribución. Para solucionar este problema es necesario aumentar el número de tallos en nuestro diagrama. Una manera sencilla de hacerlo consiste en escribir dos veces cada valor del tallo y después registrar las hojas 0, 1, 2, 3 y 4 enfrente del valor del tallo adecuado, donde aparezca por primera vez; y las hojas 5, 6, 7, 8 y 9 enfrente de este mismo valor del tallo, donde aparece la segunda vez. El diagrama doble de tallo y hojas modificado se ilustra en la tabla 1.6, donde los tallos que corresponden a las hojas 0 a 4 fueron codificados con el símbolo * y los tallos correspondientes a las hojas 5 a 9 con el símbolo •.

En cualquier problema dado debemos decidir cuáles son los valores del tallo adecuados. Esta decisión se toma hasta cierto punto de manera arbitraria, aunque debemos guiarnos por el tamaño de nuestra muestra. Por lo general elegimos entre 5 y 20 tallos. Cuanto más pequeña sea la cantidad de datos disponibles, más pequeña será nuestra elección del número de tallos. Por ejemplo, si los datos constan de números del 1 al 21,

los cuales representan el número de personas en la fila de una cafetería en 40 días laborales seleccionados al azar, y elegimos un diagrama doble de tallo y hojas, los tallos serían 0*, 0., 1*, 1. y 2*, de manera que la observación de 1 más pequeña tiene tallo 0* y hoja 1, el número 18 tiene tallo 1. y hoja 8, y la observación de 21 más grande tiene tallo 2* y hoja 1. Por otro lado, si los datos constan de números de \$18,800 a \$19,600, que representan las mejores ventas posibles de 100 automóviles nuevos, obtenidos de cierto concesionario, y elegimos un diagrama sencillo de tallo y hojas, los tallos serían 188, 189, 190, ..., 196 y las hojas contendrían ahora dos dígitos cada una. Un automóvil que se vende en \$19,385 tendría un valor de tallo de 193 y 85 en los dos dígitos de la hoja. En el diagrama de tallo y hojas, las hojas de dígitos múltiples que pertenecen al mismo tallo por lo regular están separadas por comas. En los datos generalmente se ignoran los puntos decimales cuando todos los números a la derecha del punto decimal representan hojas, como en el caso de las tablas 1.5 y 1.6. Sin embargo, si los datos constaran de números que van de 21.8 a 74.9, podríamos elegir los dígitos 2, 3, 4, 5, 6 y 7 como los tallos, de manera que un número como 48.3 tendría un valor de tallo de 4 y un valor de hoja de 8.3.

Tabla 1.6: Diagrama doble de tallo y hojas para la vida de las baterías

Tallo	Hoja	Frecuencia
1.	69	2
2*	2	1
2.	5669	4
3*	001111222333444	15
3.	5567778899	10
4*	11234	5
4.	577	3

El diagrama de tallo y hojas representa una manera eficaz de resumir los datos. Otra forma consiste en el uso de la **distribución de frecuencias**, donde los datos, agrupados en diferentes clases o intervalos, se pueden construir contando las hojas que pertenecen a cada tallo y considerando que cada tallo define un intervalo de clase. En la tabla 1.5 el tallo 1 con 2 hojas define el intervalo 1.0-1.9, que contiene 2 observaciones; el tallo 2 con 5 hojas define el intervalo 2.0-2.9, que contiene 5 observaciones; el tallo 3 con 25 hojas define el intervalo 3.0-3.9, con 25 observaciones; y el tallo 4 con 8 hojas define el intervalo 4.0-4.9, que contiene 8 observaciones. Para el diagrama doble de tallo y hojas de la tabla 1.6 los tallos definen los siete intervalos de clase 1.5-1.9, 2.0-2.4, 2.5-2.9, 3.0-3.4, 3.5-3.9, 4.0-4.4 y 4.5-4.9, con frecuencias 2, 1, 4, 15, 10, 5 y 3, respectivamente.

Histograma

Al dividir cada frecuencia de clase entre el número total de observaciones, obtenemos la proporción del conjunto de observaciones en cada una de las clases. Una tabla que lista las frecuencias relativas se denomina **distribución de frecuencias relativas**. En la tabla 1.7 se presenta la distribución de frecuencias relativas para los datos de la tabla 1.4, que muestra los puntos medios de cada intervalo de clase.

La información que brinda una distribución de frecuencias relativas en forma tabular es más fácil de entender si se presenta en forma gráfica. Con los puntos medios de

Tabla 1.7: Distribución de frecuencias relativas de la vida de las baterías

Intervalo de clase	Punto medio de la clase	Frecuencia, f	Frecuencia relativa
1.5–1.9	1.7	2	0.050
2.0–2.4	2.2	1	0.025
2.5–2.9	2.7	4	0.100
3.0–3.4	3.2	15	0.375
3.5–3.9	3.7	10	0.250
4.0–4.4	4.2	5	0.125
4.5–4.9	4.7	3	0.075

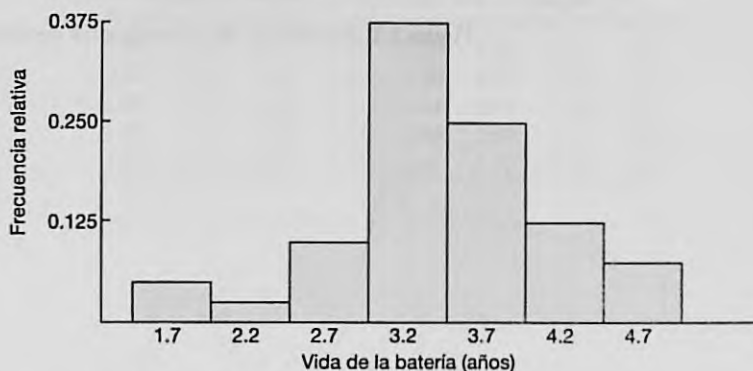


Figura 1.6: Histograma de frecuencias relativas.

cada intervalo y las frecuencias relativas correspondientes construimos un **histograma de frecuencias relativas** (figura 1.6).

Muchas distribuciones de frecuencias continuas se pueden representar gráficamente mediante la curva en forma de campana característica de la figura 1.7. Herramientas gráficas como las de las figuras 1.6 y 1.7 ayudan a comprender la naturaleza de la población. En los capítulos 5 y 6 examinaremos una propiedad de la población que se conoce como **distribución**. Aunque más adelante en este texto se proporcionará una definición más precisa de una distribución o de una **distribución de probabilidad**, aquí podemos visualizarla como la que se podría haber visto en el límite de la figura 1.7 cuando el tamaño de la muestra aumentara.

Se dice que una distribución es **simétrica** si se puede doblar a lo largo de un eje vertical de manera que ambos lados coincidan. Si una distribución carece de simetría respecto de un eje vertical, se dice que está **sesgada**. La distribución que se ilustra en la figura 1.8a se dice que está sesgada a la derecha porque tiene una cola derecha larga y una cola izquierda mucho más corta. En la figura 1.8b observamos que la distribución es simétrica; mientras que en la figura 1.8c está sesgada a la izquierda.

Al girar un diagrama de tallo y hojas en dirección contraria a la de las manecillas del reloj en un ángulo de 90° , vemos que las columnas de hojas que resultan forman una imagen parecida a un histograma. Por lo tanto, si nuestro objetivo principal al observar los datos es determinar la forma general o la forma de la distribución, rara vez será necesario construir un histograma de frecuencias relativas.

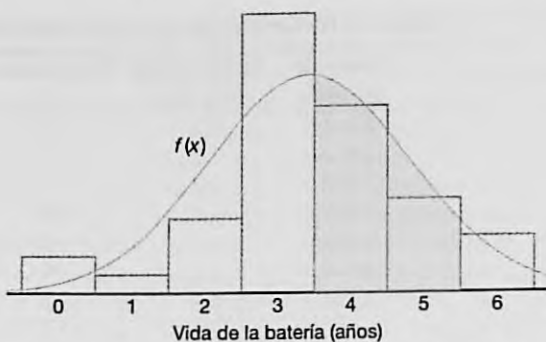


Figura 1.7: Estimación de la distribución de frecuencias.

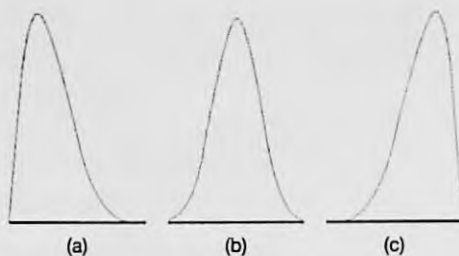


Figura 1.8: Sesgo de los datos.

Gráfica de caja y bigote o gráfica de caja

Otra presentación que es útil para reflejar propiedades de una muestra es la **gráfica de caja y bigote**, la cual encierra el *rango intercuartil* de los datos en una caja que contiene la mediana representada. El rango intercuartil tiene como extremos el percentil 75 (cuartil superior) y el percentil 25 (cuartil inferior). Además de la caja se prolongan “bigotes”, que indican las observaciones alejadas en la muestra. Para muestras razonablemente grandes la presentación indica el centro de localización, la variabilidad y el grado de asimetría.

Además, una variación denominada **gráfica de caja** puede ofrecer al observador información respecto de cuáles observaciones son **valores extremos**. Los valores extremos son observaciones que se consideran inusualmente alejadas de la masa de datos. Existen muchas pruebas estadísticas diseñadas para detectar este tipo de valores. Técnicamente se puede considerar que un valor extremo es una observación que representa un “evento raro” (existe una probabilidad pequeña de obtener un valor que esté lejos de la masa de datos). El concepto de valores extremos volverá a surgir en el capítulo 12 en el contexto del análisis de regresión.

La información visual en las gráficas de caja y bigote o en las de caja no intenta ser una prueba formal de valores extremos, más bien se considera una herramienta de diagnóstico. Aunque la determinación de cuáles observaciones son valores extremos varía de acuerdo con el tipo de software que se emplee, un procedimiento común para determinarlo consiste en utilizar un **múltiplo del rango intercuartil**. Por ejemplo, si la distancia desde la caja excede 1.5 veces el rango intercuartil (en cualquier dirección), la observación se podría considerar un valor extremo.

Ejemplo 1.5: Se midió el contenido de nicotina en una muestra aleatoria de 40 cigarrillos. Los datos se presentan en la tabla 1.8.

Tabla 1.8: Valores de nicotina para el ejemplo 1.5

1.09	1.92	2.31	1.79	2.28	1.74	1.47	1.97
0.85	1.24	1.58	2.03	1.70	2.17	2.55	2.11
1.86	1.90	1.68	1.51	1.64	0.72	1.69	1.85
1.82	1.79	2.46	1.88	2.08	1.67	1.37	1.93
1.40	1.64	2.09	1.75	1.63	2.37	1.75	1.69

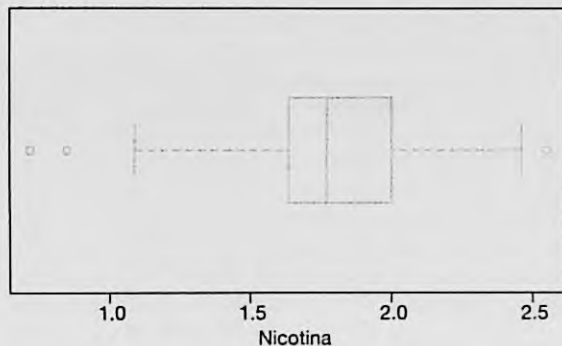


Figura 1.9: Gráfica de caja y bigote para el ejemplo 1.5.

La figura 1.9 muestra la gráfica de caja y bigote de los datos, la cual describe las observaciones 0.72 y 0.85 como valores extremos moderados en la cola inferior; en tanto que la observación 2.55 es un valor extremo moderado en la cola superior. En este ejemplo el rango intercuartil es 0.365, y 1.5 veces el rango intercuartil es 0.5475. Por otro lado, la figura 1.10 presenta un diagrama de tallo y hojas. ■

Ejemplo 1.6: Considere los datos de la tabla 1.9, que constan de 30 muestras que miden el grosor de las “asas” de latas de pintura (véase el trabajo de Hogg y Ledolter de 1992 en la bibliografía). La figura 1.11 describe una gráfica de caja y bigote para este conjunto asimétrico de datos. Observe que el bloque izquierdo es considerablemente más grande que el bloque de la derecha. La mediana es 35. El cuartil inferior es 31, mientras que el superior es 36. Advierta también que la observación alejada de la derecha está más lejos de la caja que la observación extrema de la izquierda. No hay valores extremos en este conjunto de datos. ■

El punto decimal se encuentra 1 dígito(s) a la izquierda de I

7		2
8		5
9		
10		9
11		
12		4
13		7
14		07
15		18
16		3447899
17		045599
18		2568
19		0237
20		389
21		17
22		8
23		17
24		6
25		5

Figura 1.10: Diagrama de tallo y hojas para los datos de nicotina.

Tabla 1.9: Datos para el ejemplo 1.6

Muestra	Mediciones	Muestra	Mediciones
1	29 36 39 34 34	16	35 30 35 29 37
2	29 29 28 32 31	17	40 31 38 35 31
3	34 34 39 38 37	18	35 36 30 33 32
4	35 37 33 38 41	19	35 34 35 30 36
5	30 29 31 38 29	20	35 35 31 38 36
6	34 31 37 39 36	21	32 36 36 32 36
7	30 35 33 40 36	22	36 37 32 34 34
8	28 28 31 34 30	23	29 34 33 37 35
9	32 36 38 38 35	24	36 36 35 37 37
10	35 30 37 35 31	25	36 30 35 33 31
11	35 30 35 38 35	26	35 30 29 38 35
12	38 34 35 35 31	27	35 36 30 34 36
13	34 35 33 30 34	28	35 30 36 29 35
14	40 35 34 33 35	29	38 36 35 31 31
15	34 35 38 35 30	30	30 34 40 28 30

Existen otras formas en las que las gráficas de caja y bigote, y otras presentaciones gráficas, pueden ayudar al analista. Las muestras múltiples se pueden comparar de forma gráfica. Los diagramas de los datos pueden sugerir relaciones entre las variables y las gráficas ayudan a detectar anomalías u observaciones extremas en las muestras.

Existen otros tipos diferentes de diagramas y herramientas gráficas, los cuales se estudiarán en el capítulo 8 después de presentar otros detalles teóricos.

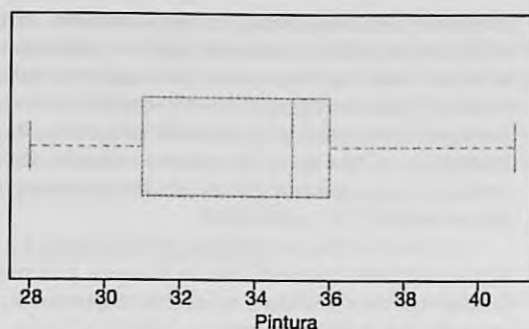


Figura 1.11: Gráfica de caja y bigote del grosor de las "asas" de latas de pintura.

Otras características distintivas de una muestra

Hay características de la distribución o de la muestra, además de las medidas del centro de localización y variabilidad, que definen aún más su naturaleza. Por ejemplo, en tanto que la mediana divide los datos (o su distribución) en dos partes, existen otras medidas que dividen partes o segmentos de la distribución que pueden ser muy útiles. Una separación en cuatro partes se hace mediante *cuartiles*, donde el tercer cuartil separa el cuarto (25%) superior del resto de los datos, el segundo cuartil es la mediana y el primer cuartil separa el cuarto (25%) inferior del resto de los datos. La distribución puede dividirse incluso más detalladamente calculando los percentiles. Tales cantidades dan al analista una noción de las denominadas *colas* de la distribución (es decir, los valores que son relativamente extremos, ya sean pequeños o grandes). Por ejemplo, el percentil 95 separa el 5% superior del 95% inferior. Para los extremos en la parte inferior o *cola inferior* de la distribución prevalecen definiciones similares. El primer percentil separa el 1% inferior del resto de la distribución. El concepto de percentiles desempeñará un papel significativo en buena parte de lo que estudiaremos en los siguientes capítulos.

1.7 Tipos generales de estudios estadísticos: diseño experimental, estudio observacional y estudio retrospectivo

En las siguientes secciones destacaremos el concepto de muestreo de una población y el uso de los métodos estadísticos para aprender o quizá para reafirmar la información relevante acerca de una población. La información que se busca y que se obtiene mediante el uso de tales métodos estadísticos a menudo influye en la toma de decisiones, así como en la resolución de problemas en diversas áreas importantes de ingeniería y científicas. Como ilustración, el ejemplo 1.3 describe un experimento sencillo, en el cual los resultados brindan ayuda para determinar los tipos de condiciones en los que no se recomienda utilizar una aleación de aluminio específica que podría ser muy vulnerable a la corrosión. Los resultados serían útiles no sólo para quienes fabrican la aleación, sino también para los clientes que consideren adquirirla. Este caso, y muchos otros que se incluyen en los capítulos 13 a 15, resaltan el concepto de condiciones experimentales diseñadas o controladas (combinaciones de condiciones de recubrimiento y humedad), que son de interés para aprender sobre algunas características o mediciones (nivel de corrosión) que

surgen de tales condiciones. En las mediciones de la corrosión se emplean métodos estadísticos que utilizan tanto medidas de tendencia central como de variabilidad. Como usted verá más adelante en este texto, tales métodos nos guían hacia un modelo estadístico como el que se examinó en la sección 1.6. En este caso el modelo se puede usar para estimar (o predecir) las medidas de la corrosión como una función de la humedad y el tipo de recubrimiento utilizado. De nuevo, para desarrollar este tipo de modelos es muy útil emplear las estadísticas descriptivas que destacan las medidas de tendencia central y de variabilidad.

La información que se ofrece en el ejemplo 1.3 ilustra de manera adecuada los tipos de preguntas de ingeniería que se plantean y se responden aplicando los métodos estadísticos que se utilizan en un diseño experimental y se presentan en este texto. Tales preguntas son las siguientes:

- i. ¿Cuál es la naturaleza del efecto de la humedad relativa sobre la corrosión de la aleación de aluminio dentro del rango de humedad relativa en este experimento?
- ii. ¿El recubrimiento químico contra la corrosión reduce los niveles de corrosión y existe alguna manera de cuantificar el efecto?
- iii. ¿Hay alguna **interacción** entre el tipo de recubrimiento y la humedad relativa que influya en la corrosión de la aleación? Si es así, ¿cómo se podría interpretar?

¿Qué es interacción?

La importancia de las preguntas i. y ii. debería quedar clara para el lector, ya que ambas tienen que ver con aspectos importantes tanto para los productores como para los usuarios de la aleación. ¿Pero qué sucede con la pregunta iii.? El concepto de *interacción* se estudiará con detalle en los capítulos 14 y 15. Considere la gráfica de la figura 1.3, la cual ejemplifica la detección de la interacción entre dos factores en un diseño experimental simple. Observe que las líneas que conectan las medias de la muestra no son paralelas. El **paralelismo** habría indicado que el efecto (visto como un resultado de la pendiente de las líneas) de la humedad relativa es igual, es decir, negativo, tanto en la condición sin recubrimiento como en la condición con recubrimiento químico contra la corrosión. Recuerde que la pendiente negativa implica que la corrosión se vuelve más pronunciada a medida que aumenta la humedad. La ausencia de paralelismo implica una interacción entre el tipo de recubrimiento y la humedad relativa. La línea casi “horizontal” para el recubrimiento contra la corrosión, opuesta a la pendiente más pronunciada para la condición sin recubrimiento, sugiere que *el recubrimiento químico contra la corrosión no sólo es benéfico (observe el desplazamiento entre las líneas), sino que la presencia del recubrimiento revela que el efecto de la humedad es despreciable*. Salta a la vista que todas estas cuestiones son muy importantes para el efecto de los dos factores individuales y para la interpretación de la interacción, si está presente.

Los modelos estadísticos son muy útiles para responder preguntas como las descritas en i, ii y iii, en donde los datos provienen de un diseño experimental. Sin embargo, no siempre se cuenta con el tiempo o los recursos que permiten usar un diseño experimental. Por ejemplo, hay muchos casos en los que las condiciones de interés para el científico o el ingeniero simplemente no se pueden implementar *debido a que es imposible controlar los factores importantes*. En el ejemplo 1.3 la humedad relativa y el tipo de recubrimiento (o la ausencia de éste) son bastante fáciles de controlar. Desde luego, se trata del rasgo distintivo de un diseño experimental. En muchos campos los factores a estudiar no pueden ser controlados por diversas razones. Un control riguroso como el del ejemplo 1.3 permite al analista confiar en que las diferencias encontradas (como en los niveles de

corrosión) se deben a los factores que se pueden controlar. Considere el ejercicio 1.6 de la página 13 como otro ejemplo. En este caso suponga que se eligen 24 especímenes de caucho de silicio y que se asignan 12 a cada uno de los niveles de temperatura de vulcanizado. Las temperaturas se controlan cuidadosamente, por lo que éste es un ejemplo de diseño experimental con **un solo factor**, que es la temperatura de vulcanizado. Se podría suponer que las diferencias encontradas en la media de la resistencia a la tensión son atribuibles a las diferentes temperaturas de vulcanizado.

¿Qué sucede si no se controlan los factores?

Suponga que los factores no se controlan y que *no hay asignación aleatoria* a los tratamientos específicos para las unidades experimentales, y que se necesita obtener información a partir de un conjunto de datos. Como ejemplo considere un estudio donde el interés se centra en la relación entre los niveles de colesterol sanguíneo y la cantidad de sodio medida en la sangre. Durante cierto periodo se revisó el colesterol sanguíneo y el nivel de sodio de un grupo de individuos. En efecto, es posible obtener alguna información útil de tal conjunto de datos. Sin embargo, debería quedar claro que no es posible hacer un control estricto de los niveles de sodio. De manera ideal, los sujetos deberían dividirse aleatoriamente en dos grupos, donde uno fuera el asignado a un nivel alto específico de sodio en la sangre, y el otro a un nivel bajo específico de sodio en la sangre, pero es obvio que esto no es posible. Evidentemente los cambios en los niveles de colesterol se deben a cambios en uno o diversos factores que no se controlaron. Este tipo de estudio, sin control de factores, se denomina **estudio observacional**, el cual la mayoría de las veces implica una situación en que los sujetos se observan a través del tiempo.

Los estudios biológicos y biomédicos a menudo tienen que ser observacionales. Sin embargo, este tipo de estudios no se restringen a dichas áreas. Por ejemplo, considere un estudio diseñado para determinar la influencia de la temperatura ambiental sobre la energía eléctrica que consumen las instalaciones de una planta química. Es evidente que los niveles de la temperatura ambiental no se pueden controlar, por lo tanto, la única manera en que se puede supervisar la estructura de los datos es a partir de los datos de la planta a través del tiempo.

Es importante destacar que una diferencia básica entre un experimento bien diseñado y un estudio observacional es la dificultad para determinar la causa y el efecto verdaderos con este último. Asimismo, las diferencias encontradas en la reacción fundamental (por ejemplo, niveles de corrosión, colesterol sanguíneo, consumo de energía eléctrica en una planta) podrían deberse a otros factores subyacentes que no se controlaron. De manera ideal, en un diseño experimental los *factores perturbadores* serían compensados mediante el proceso de aleatoriedad. En realidad, los cambios en los niveles de colesterol sanguíneo podrían deberse a la ingestión de grasa, a la realización de actividad física, etc. El consumo de energía eléctrica podría estar afectado por la cantidad de bienes producidos o incluso por la pureza de éstos.

Otra desventaja de los estudios observacionales, que a menudo se ignora cuando éstos se comparan con experimentos cuidadosamente diseñados, es que, a diferencia de estos últimos, los observacionales están a merced de circunstancias no controladas, naturales, ambientales o de otros tipos, que repercuten en los niveles de los factores de interés. Por ejemplo, en el estudio biomédico acerca de la influencia de los niveles de sodio en la sangre sobre el colesterol sanguíneo es posible que haya, de hecho, una influencia significativa, pero el conjunto de datos específico usado no involucró la suficiente variación observada en los niveles de sodio debido a la naturaleza de los sujetos elegidos. Desde luego, en un diseño experimental el analista elige y controla los niveles de los factores.

Un tercer tipo de estudio estadístico que podría ser muy útil, pero que tiene notables desventajas cuando se le compara con un diseño experimental, es el **estudio retrospectivo**. Esta clase de estudio emplea estrictamente **datos históricos**, que se obtienen durante un periodo específico. Una ventaja evidente de los datos retrospectivos es el bajo costo de la recopilación de datos. Sin embargo, como se podría esperar, también tiene desventajas claras:

- i. La validez y la confiabilidad de los datos históricos a menudo son cuestionables.
- ii. Si el tiempo es un aspecto relevante en la estructura de los datos, podría haber datos faltantes.
- iii. Podrían existir errores en la recopilación de los datos que no se conocen.
- iv. De nuevo, como en el caso de los datos observacionales, no hay control en los rangos de las variables a medir (es decir, en los factores a estudiar). De hecho, las variaciones que se encuentran en los datos históricos a menudo no son significativas para estudios actuales.

En la sección 1.6 se puso cierto énfasis en los modelos de las relaciones entre variables. Presentamos el concepto de análisis de regresión, el cual se estudia en los capítulos 11 y 12, y se considera como una forma del análisis de datos para los diseños experimentales que se examinarán en los capítulos 14 y 15. En la sección 1.6 se utilizó a modo de ejemplo un modelo que relaciona la media poblacional de la resistencia a la tensión de la tela con los porcentajes de algodón, en el cual 20 especímenes de tela representaban las unidades experimentales. En este caso, los datos provienen de un diseño experimental simple, en el que los porcentajes de algodón individuales fueron seleccionados por el científico.

Con frecuencia, tanto los datos observacionales como los retrospectivos se utilizan para observar las relaciones entre variables a través de los procedimientos de construcción de modelos que se estudiarán en los capítulos 11 y 12. Aunque las ventajas de los diseños experimentales se pueden aplicar cuando la finalidad es la construcción de un modelo estadístico, hay muchas áreas en las que no es posible diseñar experimentos, de manera que *habrá que utilizar los datos históricos u observacionales*. Aquí nos referimos al conjunto de datos históricos que se incluye en el ejercicio 12.5 de la página 450. El objetivo es construir un modelo que dé como resultado una ecuación o relación que vincule el consumo mensual de energía eléctrica con la temperatura ambiental promedio, x_1 , el número de días en el mes, x_2 , la pureza promedio del producto, x_3 y las toneladas de bienes producidos, x_4 . Se trata de los datos históricos del año anterior.

Ejercicios

1.13 Un fabricante de componentes electrónicos se interesa en determinar el tiempo de vida de cierto tipo de batería. Una muestra, en horas de vida, es como la siguiente:

123, 116, 122, 110, 175, 126, 125, 111, 118, 117.

- a) Calcule la media y la mediana de la muestra.
- b) ¿Qué característica en este conjunto de datos es la responsable de la diferencia sustancial entre ambas?

1.14 Un fabricante de neumáticos quiere determinar el diámetro interior de un neumático de cierto grado de calidad. Idealmente el diámetro sería de 570 mm. Los datos son los siguientes:

572, 572, 573, 568, 569, 575, 565, 570.

- a) Calcule la media y la mediana de la muestra.
- b) Obtenga la varianza, la desviación estándar y el rango de la muestra.
- c) Con base en los estadísticos calculados en los incisos a) y b), ¿qué comentaría acerca de la calidad de los neumáticos?

1.15 Cinco lanzamientos independientes de una moneda tienen como resultado *cinco caras*. Resulta que si la moneda es legal, la probabilidad de este resultado es $(1/2)^5 = 0.03125$. ¿Proporciona esto evidencia sólida

de que la moneda no es legal? Comente y utilice el concepto de valor- P que se analizó en la sección 1.1.

1.16 Muestre que las n piezas de información en $\sum_{i=1}^n (x_i - \bar{x})^2$ no son independientes; es decir, demuestre que

$$\sum_{i=1}^n (x_i - \bar{x}) = 0.$$

1.17 Se realiza un estudio acerca de los efectos del tabaquismo sobre los patrones de sueño. La medición que se observa es el tiempo, en minutos, que toma quedar dormido. Se obtienen los siguientes datos:

Fumadores:	69.3	56.0	22.1	47.6
	53.2	48.1	52.7	34.4
	60.2	43.8	23.2	13.8
No fumadores:	28.6	25.1	26.4	34.9
	29.8	28.4	38.5	30.2
	30.6	31.8	41.6	21.1
	36.0	37.9	13.9	

- Calcule la media de la muestra para cada grupo.
- Calcule la desviación estándar de la muestra para cada grupo.
- Elabore una gráfica de puntos de los conjuntos de datos A y B en la misma línea.
- Comente qué clase de efecto parece tener el hecho de fumar sobre el tiempo que se requiere para quedarse dormido.

1.18 Las siguientes puntuaciones representan la calificación en el examen final para un curso de estadística elemental:

23	60	79	32	57	74	52	70	82
36	80	77	81	95	41	65	92	85
55	76	52	10	64	75	78	25	80
98	81	67	41	71	83	54	64	72
88	62	74	43	60	78	89	76	84
48	84	90	15	79	34	67	17	82
69	74	63	80	85	61			

- Elabore un diagrama de tallo y hojas para las calificaciones del examen, donde los tallos sean 1, 2, 3, ..., 9.
- Elabore un histograma de frecuencias relativas, trace un estimado de la gráfica de la distribución y analice la asimetría de la distribución.
- Calcule la media, la mediana y la desviación estándar de la muestra.

1.19 Los siguientes datos representan la duración de vida, en años, medida al entero más cercano, de 30 bombas de combustible similares.

2.0	3.0	0.3	3.3	1.3	0.4
0.2	6.0	5.5	6.5	0.2	2.3
1.5	4.0	5.9	1.8	4.7	0.7

4.5	0.3	1.5	0.5	2.5	5.0
1.0	6.0	5.6	6.0	1.2	0.2

- Construya un diagrama de tallo y hojas para la vida, en años, de las bombas de combustible, utilizando el dígito a la izquierda del punto decimal como el tallo para cada observación.
- Determine una distribución de frecuencias relativas.
- Calcule la media, el rango y la desviación estándar de la muestra.

1.20 Los siguientes datos representan la duración de la vida, en segundos, de 50 moscas de la fruta que se someten a un nuevo aerosol en un experimento de laboratorio controlado.

17	20	10	9	23	13	12	19	18	24
12	14	6	9	13	6	7	10	13	7
16	18	8	3	3	32	9	7	10	11
13	7	18	7	10	4	27	19	16	8
7	10	5	14	15	10	9	6	7	15

- Elabore un diagrama doble de tallo y hojas para el periodo de vida de las moscas de la fruta usando los tallos 0★, 0•, 1★, 1•, 2★, 2• y 3★ de manera que los tallos codificados con los símbolos ★ y • se asocien, respectivamente, con las hojas 0 a 4 y 5 a 9.
- Determine una distribución de frecuencias relativas.
- Construya un histograma de frecuencias relativas.
- Calcule la mediana.

1.21 La duración de fallas eléctricas, en minutos, se presenta en la siguiente tabla.

22	18	135	15	90	78	69	98	102
83	55	28	121	120	13	22	124	112
70	66	74	89	103	24	21	112	21
40	98	87	132	115	21	28	43	37
50	96	118	158	74	78	83	93	95

- Calcule la media y la mediana muestrales de las duraciones de la falla eléctrica.
- Calcule la desviación estándar de las duraciones de la falla eléctrica.

1.22 Los siguientes datos son las mediciones del diámetro de 36 cabezas de remache en centésimos de una pulgada.

6.72	6.77	6.82	6.70	6.78	6.70	6.62	6.75
6.66	6.66	6.64	6.76	6.73	6.80	6.72	6.76
6.76	6.68	6.66	6.62	6.72	6.76	6.70	6.78
6.76	6.67	6.70	6.72	6.74	6.81	6.79	6.78
6.66	6.76	6.76	6.72				

- Calcule la media y la desviación estándar de la muestra.
- Construya un histograma de frecuencias relativas para los datos.

- c) Comente si existe o no una indicación clara de que la muestra proviene de una población que tiene una distribución en forma de campana.

1.23 En 20 automóviles elegidos aleatoriamente, se tomaron las emisiones de hidrocarburos en velocidad en vacío, en partes por millón (ppm), para modelos de 1980 y 1990.

Modelos 1980:

141 359 247 940 882 494 306 210 105 880
200 223 188 940 241 190 300 435 241 380

Modelos 1990:

140 160 20 20 223 60 20 95 360 70
220 400 217 58 235 380 200 175 85 65

- a) Construya una gráfica de puntos como la de la figura 1.1.
b) Calcule la media de la muestra para los dos años y sobreponga las dos medias en las gráficas.
c) Comente sobre lo que indica la gráfica de puntos respecto de si cambiaron o no las emisiones poblacionales de 1980 a 1990. Utilice el concepto de variabilidad en sus comentarios.

1.24 Los siguientes son datos históricos de los sueldos del personal (dólares por alumno) en 30 escuelas seleccionadas de la región este de Estados Unidos a principios de la década de 1970.

3.79 2.99 2.77 2.91 3.10 1.84 2.52 3.22
2.45 2.14 2.67 2.52 2.71 2.75 3.57 3.85
3.36 2.05 2.89 2.83 3.13 2.44 2.10 3.71
3.14 3.54 2.37 2.68 3.51 3.37

- a) Calcule la media y la desviación estándar de la muestra.
b) Utilice los datos para elaborar un histograma de frecuencias relativas.
c) Construya un diagrama de tallo y hojas con los datos.

1.25 El siguiente conjunto de datos se relaciona con el ejercicio 1.24 y representa el porcentaje de las familias que se ubican en el nivel superior de ingresos en las mismas escuelas individuales y con el mismo orden del ejercicio 1.24.

72.2 31.9 26.5 29.1 27.3 8.6 22.3 26.5
20.4 12.8 25.1 19.2 24.1 58.2 68.1 89.2
55.1 9.4 14.5 13.9 20.7 17.9 8.5 55.4
38.1 54.2 21.5 26.2 59.1 43.3

- a) Calcule la media de la muestra.
b) Calcule la mediana de la muestra.
c) Construya un histograma de frecuencias relativas con los datos.
d) Determine la media recortada al 10%. Compárela con los resultados de los incisos a) y b) y exprese su comentario.

1.26 Suponga que le interesa emplear los conjuntos de datos de los ejercicios 1.24 y 1.25 para derivar un modelo

que prediga los salarios del personal como una función del porcentaje de familias en un nivel alto de ingresos para los sistemas escolares actuales. Comente sobre cualquier desventaja de llevar a cabo este tipo de análisis.

1.27 Se realizó un estudio para determinar la influencia del desgaste, y , de un cojinete como una función de la carga, x , sobre el cojinete. Para este estudio se utilizó un diseño experimental con tres niveles de carga: 700 lb, 1000 lb y 1300 lb. En cada nivel se utilizaron cuatro especímenes y las medias muestrales fueron 210, 325 y 375, respectivamente.

- a) Grafique el promedio de desgaste contra la carga.
b) A partir de la gráfica del inciso a), ¿consideraría que hay una relación entre desgaste y carga?
c) Suponga que tenemos los siguientes valores individuales de desgaste para cada uno de los cuatro especímenes en los respectivos niveles de carga. (Vea los datos que siguen). Grafique los resultados de desgaste para todos los especímenes contra los tres valores de carga.
d) A partir de la gráfica del inciso c), ¿consideraría que hay una relación clara? Si su respuesta difiere de la del inciso b), explique por qué.

	x		
	700	1000	1300
y_1	145	250	150
y_2	105	195	180
y_3	260	375	420
y_4	330	480	750
	$\bar{y}_1 = 210$	$\bar{y}_2 = 325$	$\bar{y}_3 = 375$

1.28 En Estados Unidos y otros países muchas compañías de manufactura utilizan partes moldeadas como componentes de un proceso. La contracción a menudo es un problema importante. Por consiguiente, un dado de metal moldeado para una parte se construye más grande que el tamaño nominal con el fin de permitir su contracción. En un estudio de moldeado por inyección se descubrió que en la contracción influyen múltiples factores, entre los cuales están la velocidad de la inyección en pies/segundo y la temperatura de moldeado en °C. Los dos conjuntos de datos siguientes muestran los resultados del diseño experimental, en donde la velocidad de inyección se mantuvo a dos niveles (bajo y alto) y la temperatura de moldeado se mantuvo constante en un nivel bajo. La contracción se midió en $\text{cm} \times 10^4$. Los valores de contracción a una velocidad de inyección baja fueron:

72.68 72.62 72.58 72.48 73.07
72.55 72.42 72.84 72.58 72.92

Los valores de contracción a una velocidad de inyección alta fueron:

71.62 71.68 71.74 71.48 71.55
71.52 71.71 71.56 71.70 71.50

- a) Construya una gráfica de puntos para ambos conjuntos de datos en la misma gráfica. Sobre ésta indique ambas medias de la contracción, tanto para la velocidad de inyección baja como para la velocidad de inyección alta.
- b) Con base en los resultados de la gráfica del inciso a), y considerando la ubicación de las dos medias y su sentido de variabilidad, ¿cuál es su conclusión respecto del efecto de la velocidad de inyección sobre la contracción a una temperatura de moldeo baja?

1.29 Utilice los datos del ejercicio 1.24 para elaborar una gráfica de caja.

1.30 A continuación se presentan los tiempos de vida, en horas, de 50 lámparas incandescentes, con esmerilado interno, de 40 watts y 110 voltios, los cuales se tomaron de pruebas forzadas de vida:

919	1196	785	1126	936	918
1156	920	948	1067	1092	1162
1170	929	950	905	972	1035
1045	855	1195	1195	1340	1122
938	970	1237	956	1102	1157
978	832	1009	1157	1151	1009
765	958	902	1022	1333	811
1217	1085	896	958	1311	1037
702	923				

Elabore una gráfica de puntos para estos datos.

1.31 Considere la situación del ejercicio 1.28, pero ahora utilice el siguiente conjunto de datos, en el cual la contracción se mide de nuevo a una velocidad de inyección baja y a una velocidad de inyección alta. Sin embargo, esta vez la temperatura de moldeo se aumenta a un nivel "alto" y se mantiene constante.

Los valores de la contracción a una velocidad de inyección baja fueron:

76.20	76.09	75.98	76.15	76.17
75.94	76.12	76.18	76.25	75.82

Los valores de la contracción a una velocidad de inyección alta fueron:

93.25	93.19	92.87	93.29	93.37
92.98	93.47	93.75	93.89	91.62

- a) Igual que en el ejercicio 1.28, elabore una gráfica de puntos con ambos conjuntos de datos en la misma gráfica e identifique las dos medias (es decir, la contracción media para la velocidad de inyección baja y para la velocidad de inyección alta).
- b) Igual que en el ejercicio 1.28, comente sobre la influencia de la velocidad de inyección en la contracción para la temperatura de moldeo alta. Tome en cuenta la posición de las dos medias y la variabilidad de cada media.
- c) Compare su conclusión en el inciso b) actual con la del inciso b) del ejercicio 1.28, en el cual la temperatura de moldeo se mantuvo a un nivel bajo. ¿Diría que hay interacción entre la velocidad de inyección y la temperatura de moldeo? Explique su respuesta.

1.32 Utilice los resultados de los ejercicios 1.28 y 1.31 para crear una gráfica que ilustre la interacción evidente entre los datos. Use como guía la gráfica de la figura 1.3 del ejemplo 1.3. ¿El tipo de información que se encontró en los ejercicios 1.28 y 1.31 se habría encontrado en un estudio observacional en el que el analista no hubiera tenido control sobre la velocidad de inyección ni sobre la temperatura de moldeo? Explique su respuesta.

1.33 Proyecto de grupo: Registre el tamaño de calzado que usa cada estudiante de su grupo. Utilice las medias y las varianzas muestrales, así como los tipos de gráficas que se estudiaron en este capítulo, para resumir cualquier característica que revele una diferencia entre las distribuciones del tamaño del calzado de hombres y mujeres. Haga lo mismo con la estatura de cada estudiante de su grupo.

Capítulo 2

Probabilidad

2.1 Espacio muestral

En el estudio de la estadística tratamos básicamente con la presentación e interpretación de **resultados fortuitos** que ocurren en un estudio planeado o en una investigación científica. Por ejemplo, en Estados Unidos, y con la finalidad de justificar la instalación de un semáforo, se podría registrar el número de accidentes que ocurren mensualmente en la intersección de Driftwood Lane y Royal Oak Drive; en una fábrica se podrían clasificar los artículos que salen de la línea de ensamble como “defectuosos” o “no defectuosos”; en una reacción química se podría revisar el volumen de gas que se libera cuando se varía la concentración de un ácido. Por ello, quienes se dedican a la estadística a menudo manejan datos numéricos que representan conteos o mediciones, o **datos categóricos** que se podrían clasificar de acuerdo con algún criterio.

En este capítulo, al referirnos a cualquier registro de información, ya sea numérico o categórico, utilizaremos el término **observación**. Por consiguiente, los números 2, 0, 1 y 2, que representan el número de accidentes que ocurrieron cada mes, de enero a abril, durante el año pasado en la intersección de Driftwood Lane y Royal Oak Drive, constituyen un conjunto de observaciones. Lo mismo ocurre con los datos categóricos N, D, N, N y D , que representan los artículos defectuosos o no defectuosos cuando se inspeccionan cinco artículos y se registran como observaciones.

Los estadísticos utilizan la palabra **experimento** para describir cualquier proceso que genere un conjunto de datos. Un ejemplo simple de experimento estadístico es el lanzamiento de una moneda al aire. En tal experimento sólo hay dos resultados posibles: cara o cruz. Otro experimento podría ser el lanzamiento de un misil y la observación de la velocidad a la que se desplaza en tiempos específicos. Las opiniones de los votantes respecto de un nuevo impuesto sobre las ventas también se pueden considerar como observaciones de un experimento. En estadística nos interesan, en particular, las observaciones que se obtienen al repetir varias veces un experimento. En la mayoría de los casos los resultados dependerán del azar, por lo tanto, no se pueden predecir con certeza. Si un químico realizara un análisis varias veces en las mismas condiciones, obtendría diferentes medidas, las cuales indicarían un elemento de probabilidad en el procedimiento experimental. Aun cuando lancemos una moneda al aire repetidas veces, no podemos tener la certeza de que en un lanzamiento determinado obtendremos cara como resultado. Sin embargo, conocemos el conjunto completo de posibilidades para cada lanzamiento.

Dado lo expuesto en la sección 1.7, en la que se revisaron tres tipos de estudios estadísticos y se dieron varios ejemplos de cada uno, ya deberíamos estar familiarizados con el alcance del término experimento. En cada uno de los tres casos, *diseños experimentales*, *estudios observacionales* y *estudios retrospectivos*, el resultado final fue un conjunto

de datos que, por supuesto, está sujeto a la **incertidumbre**. Aunque sólo uno de ellos tiene la palabra *experimento* en su descripción, el proceso de generar los datos o el proceso de observarlos forma parte de un experimento. El estudio de la corrosión expuesto en la sección 1.2 ciertamente implica un experimento en el que los datos son representados por las mediciones de la corrosión. El ejemplo de la sección 1.7, en el que se observó el colesterol y el sodio en la sangre de un conjunto de individuos, representó un estudio observacional (como lo opuesto a un *diseño* experimental) en el que el proceso incluso generó datos y un resultado sujeto a la incertidumbre; por lo tanto, se trata de un experimento. Un tercer ejemplo en la sección 1.7 consistió en un estudio retrospectivo, en el cual se observaron datos históricos sobre el consumo de energía eléctrica por mes y el promedio mensual de la temperatura ambiental. Aun cuando los datos pueden haber estado archivados durante décadas, el proceso se seguirá considerando un experimento.

Definición 2.1: Al conjunto de todos los resultados posibles de un experimento estadístico se le llama **espacio muestral** y se representa con el símbolo S .

A cada resultado en un espacio muestral se le llama **elemento** o **miembro** del espacio muestral, o simplemente **punto muestral**. Si el espacio muestral tiene un número finito de elementos, podemos *listar* los miembros separados por comas y encerrarlos entre llaves. Por consiguiente, el espacio muestral S , de los resultados posibles cuando se lanza una moneda al aire, se puede escribir como

$$S = \{H, T\},$$

en donde H y T corresponden a “caras” y “cruces”, respectivamente.

Ejemplo 2.1: Considere el experimento de lanzar un dado. Si nos interesara el número que aparece en la cara superior, el espacio muestral sería

$$S_1 = \{1, 2, 3, 4, 5, 6\}$$

Si sólo estuviéramos interesados en si el número es par o impar, el espacio muestral sería simplemente

$$S_2 = \{\text{par, impar}\} \quad \blacksquare$$

El ejemplo 2.1 ilustra el hecho de que se puede usar más de un espacio muestral para describir los resultados de un experimento. En este caso, S_1 brinda más información que S_2 . Si sabemos cuál elemento ocurre en S_1 , podremos indicar cuál resultado tiene lugar en S_2 ; no obstante, saber lo que pasa en S_2 no ayuda mucho a determinar qué elemento ocurre en S_1 . En general, lo deseable sería utilizar un espacio muestral que proporcione la mayor información acerca de los resultados del experimento. En algunos experimentos es útil listar los elementos del espacio muestral de forma sistemática utilizando un **diagrama de árbol**.

Ejemplo 2.2: Un experimento consiste en lanzar una moneda y después lanzarla una segunda vez si sale cara. Si en el primer lanzamiento sale cruz, entonces se lanza un dado una vez. Para listar los elementos del espacio muestral que proporciona la mayor información construimos el diagrama de árbol de la figura 2.1. Las diversas trayectorias a lo largo de las ramas del árbol dan los distintos puntos muestrales. Si empezamos con la rama superior izquierda y nos movemos a la derecha a lo largo de la primera trayectoria, obtenemos el punto muestral HH , que indica la posibilidad de que ocurran caras en dos lanzamientos sucesivos de la moneda. De igual manera, el punto muestral $T3$ indica la posibilidad de que la moneda muestre una cruz seguida por un 3 en el lanzamiento del dado. Al seguir todas las trayectorias, vemos que el espacio muestral es

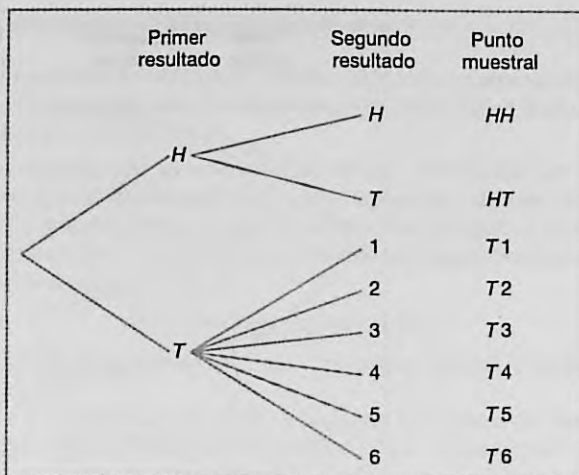


Figura 2.1: Diagrama de árbol para el ejemplo 2.2.

$$S = \{HH, HT, T1, T2, T3, T4, T5, T6\}.$$

Muchos de los conceptos de este capítulo se ilustran mejor con ejemplos que involucren el uso de dados y cartas. Es particularmente importante utilizar estas aplicaciones al comenzar el proceso de aprendizaje, ya que facilitan el flujo de esos conceptos nuevos en ejemplos científicos y de ingeniería como el siguiente.

Ejemplo 2.3: Suponga que se seleccionan, de forma aleatoria, tres artículos de un proceso de fabricación. Cada artículo se inspecciona y se clasifica como defectuoso, D , o no defectuoso, N . Para listar los elementos del espacio muestral que brinde la mayor información, construimos el diagrama de árbol de la figura 2.2, de manera que las diversas trayectorias a lo largo de las ramas del árbol dan los distintos puntos muestrales. Al comenzar con la primera trayectoria, obtenemos el punto muestral DDD , que indica la posibilidad de que los tres artículos inspeccionados estén defectuosos. Conforme continuamos a lo largo de las demás trayectorias, vemos que el espacio muestral es

$$S = \{DDD, DDN, DND, DNN, NDD, NDN, NND, NNN\}.$$

Los espacios muestrales con un número grande o infinito de puntos muestrales se describen mejor mediante un **enunciado** o **método de la regla**. Por ejemplo, si el conjunto de resultados posibles de un experimento fuera el conjunto de ciudades en el mundo con una población de más de un millón de habitantes, nuestro espacio muestral se escribiría como

$$S = \{x \mid x \text{ es una ciudad con una población de más de un millón de habitantes}\},$$

que se lee “ S es el conjunto de todas las x , tales que x es una ciudad con una población de más de un millón de habitantes”. La barra vertical se lee como “tal que”. De manera similar, si S es el conjunto de todos los puntos (x, y) sobre los límites o el interior de un círculo de radio 2 con centro en el origen, escribimos la **regla**

$$S = \{(x, y) \mid x^2 + y^2 \leq 4\}.$$

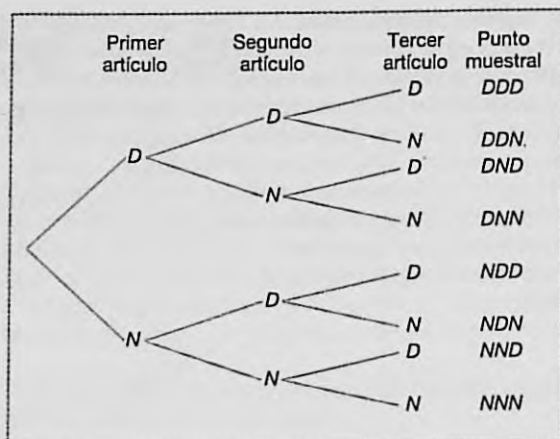


Figura 2.2: Diagrama de árbol para el ejemplo 2.3.

Nuestra elección respecto a describir el espacio muestral utilizando el método de la regla o listando los elementos dependerá del problema específico en cuestión. El método de la regla tiene ventajas prácticas, sobre todo en el caso de muchos experimentos en los que listar se vuelve una tarea tediosa.

Considere la situación del ejemplo 2.3, en el que los artículos que salen del proceso de fabricación están defectuosos, D , o no defectuosos, N . Hay muchos procedimientos estadísticos importantes llamados planes de muestreo, que determinan si un "lote" de artículos se considera o no satisfactorio. Este tipo de planes implican tomar muestras hasta obtener k artículos defectuosos. Suponga que el experimento consiste en tomar muestras de artículos, de forma aleatoria, hasta que salga uno defectuoso. En este caso el espacio muestral sería

$$S = \{D, ND, NND, NNND, \dots\}.$$

2.2 Eventos

En cualquier experimento dado, podríamos estar interesados en la ocurrencia de ciertos **eventos**, más que en la ocurrencia de un elemento específico en el espacio muestral. Por ejemplo, quizás estemos interesados en el evento A , en el cual el resultado de lanzar un dado es divisible entre 3. Esto ocurrirá si el resultado es un elemento del subconjunto $A = \{3, 6\}$ del espacio muestral S_1 del ejemplo 2.1. Otro ejemplo: podríamos estar interesados en el evento B de que el número de artículos defectuosos sea mayor que 1 en el ejemplo 2.3. Esto ocurrirá si el resultado es un elemento del subconjunto

$$B = \{DDN, DND, NDD, DDD\}$$

del espacio muestral S .

Para cada evento asignamos un conjunto de puntos muestrales, que constituye un subconjunto del espacio muestral. Este subconjunto representa la totalidad de los elementos para los que el evento es cierto.

Definición 2.2: Un evento es un subconjunto de un espacio muestral.

Ejemplo 2.4: Dado el espacio muestral $S = \{t \mid t \geq 0\}$, donde t es la vida en años de cierto componente electrónico, el evento A de que el componente falle antes de que finalice el quinto año es el subconjunto $A = \{t \mid 0 \leq t < 5\}$.

Es posible concebir que un evento puede ser un subconjunto que incluye todo el espacio muestral S , o un subconjunto de S que se denomina **conjunto vacío** y se denota con el símbolo ϕ , que no contiene ningún elemento. Por ejemplo, si en un experimento biológico permitimos que A sea el evento de detectar un organismo microscópico a simple vista, entonces $A = \phi$. También, si

$$B = \{x \mid x \text{ es un factor par de } 7\},$$

entonces B debe ser el conjunto vacío, pues los únicos factores posibles de 7 son los números nones 1 y 7.

Considere un experimento en el que se registran los hábitos de tabaquismo de los empleados de una empresa industrial. Un posible espacio muestral podría clasificar a un individuo como no fumador, fumador ocasional, fumador moderado o fumador empedernido. Si se determina que el subconjunto de los fumadores sea un evento, entonces la totalidad de los no fumadores corresponderá a un evento diferente, también subconjunto de S , que se denomina **complemento** del conjunto de fumadores.

Definición 2.3: El **complemento** de un evento A respecto de S es el subconjunto de todos los elementos de S que no están en A . Denotamos el complemento de A mediante el símbolo A' .

Ejemplo 2.5: Sea R el evento de que se seleccione una carta roja de una baraja ordinaria de 52 cartas, y sea S toda la baraja. Entonces R' es el evento de que la carta seleccionada de la baraja no sea una roja sino una negra.

Ejemplo 2.6: Considere el espacio muestral

$$S = \{\text{libro, teléfono celular, mp3, papel, papelería, computadora}\}.$$

Sea $A = \{\text{libro, papelería, computadora, papel}\}$. Entonces, el complemento de A es $A' = \{\text{teléfono celular, mp3}\}$.

Consideremos ahora ciertas operaciones con eventos que darán como resultado la formación de nuevos eventos. Estos eventos nuevos serán subconjuntos del mismo espacio muestral que los eventos dados. Suponga que A y B son dos eventos que se asocian con un experimento. En otras palabras, A y B son subconjuntos del mismo espacio muestral S . Por ejemplo, en el lanzamiento de un dado podríamos hacer que A sea el evento de que ocurra un número par y B el evento de que aparezca un número mayor que 3. Entonces, los subconjuntos $A = \{2, 4, 6\}$ y $B = \{4, 5, 6\}$ son subconjuntos del mismo espacio muestral

$$S = \{1, 2, 3, 4, 5, 6\}.$$

Observe que *tanto* A como B ocurrirán en un lanzamiento dado si el resultado es un elemento del subconjunto $\{4, 6\}$, el cual es precisamente la **intersección** de A y B .

Definición 2.4: La **intersección** de dos eventos A y B , que se denota con el símbolo $A \cap B$, es el evento que contiene todos los elementos que son comunes a A y a B .

Ejemplo 2.7: Sea E el evento de que una persona seleccionada al azar en un salón de clases sea estudiante de ingeniería, y sea F el evento de que la persona sea mujer. Entonces $E \cap F$ es el evento de todas las estudiantes mujeres de ingeniería en el salón de clases.

Ejemplo 2.8: Sean $V = \{a, e, i, o, u\}$ y $C = \{l, r, s, t\}$; entonces, se deduce que $V \cap C = \phi$. Es decir, V y C no tienen elementos comunes, por lo tanto, no pueden ocurrir de forma simultánea. ─

Para ciertos experimentos estadísticos no es nada extraño definir dos eventos, A y B , que no pueden ocurrir de forma simultánea. Se dice entonces que los eventos A y B son **mutuamente excluyentes**. Expresado de manera más formal, tenemos la siguiente definición:

Definición 2.5: Dos eventos A y B son **mutuamente excluyentes** o **disjuntos** si $A \cap B = \phi$; es decir, si A y B no tienen elementos en común.

Ejemplo 2.9: Una empresa de televisión por cable ofrece programas en ocho diferentes canales, tres de los cuales están afiliados con ABC, dos con NBC y uno con CBS. Los otros dos son un canal educativo y el canal de deportes ESPN. Suponga que un individuo que se suscribe a este servicio enciende un televisor sin seleccionar de antemano el canal. Sea A el evento de que el programa pertenezca a la cadena NBC y B el evento de que pertenezca a la cadena CBS. Como un programa de televisión no puede pertenecer a más de una cadena, los eventos A y B no tienen programas en común. Por lo tanto, la intersección $A \cap B$ no contiene programa alguno y, en consecuencia, los eventos A y B son mutuamente excluyentes. ─

A menudo nos interesamos en la ocurrencia de al menos uno de dos eventos asociados con un experimento. Por consiguiente, en el experimento del lanzamiento de un dado, si

$$A = \{2, 4, 6\} \text{ y } B = \{4, 5, 6\},$$

podríamos estar interesados en que ocurran A o B , o en que ocurran tanto A como B . Tal evento, que se llama **unión** de A y B , ocurrirá si el resultado es un elemento del subconjunto $\{2, 4, 5, 6\}$.

Definición 2.6: La **unión** de dos eventos A y B , que se denota con el símbolo $A \cup B$, es el evento que contiene todos los elementos que pertenecen a A o a B , o a ambos.

Ejemplo 2.10: Sea $A = \{a, b, c\}$ y $B = \{b, c, d, e\}$; entonces, $A \cup B = \{a, b, c, d, e\}$. ─

Ejemplo 2.11: Sea P el evento de que un empleado de una empresa petrolera seleccionado al azar fume cigarrillos. Sea Q el evento de que el empleado seleccionado ingiera bebidas alcohólicas. Entonces, el evento $P \cup Q$ es el conjunto de todos los empleados que beben o fuman, o que hacen ambas cosas. ─

Ejemplo 2.12: Si $M = \{x \mid 3 < x < 9\}$ y $N = \{y \mid 5 < y < 12\}$, entonces,

$$M \cup N = \{z \mid 3 < z < 12\}. \quad \text{─}$$

La relación entre eventos y el correspondiente espacio muestral se puede ilustrar de forma gráfica utilizando **diagramas de Venn**. En un diagrama de Venn representamos el espacio muestral como un rectángulo y los eventos con círculos trazados dentro del rectángulo. De esta forma, en la figura 2.3 vemos que

$$\begin{aligned} A \cap B &= \text{regiones 1 y 2,} \\ B \cap C &= \text{regiones 1 y 3,} \end{aligned}$$

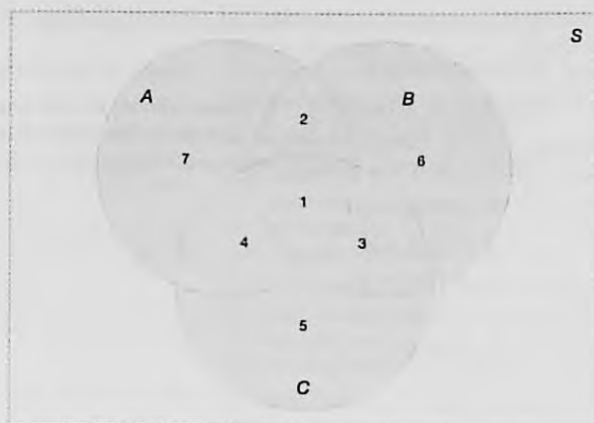


Figura 2.3: Eventos representados por varias regiones.

$$\begin{aligned}
 A \cup C &= \text{regiones } 1, 2, 3, 4, 5 \text{ y } 7, \\
 B' \cap A &= \text{regiones } 4 \text{ y } 7, \\
 A \cap B \cap C &= \text{región } 1, \\
 (A \cup B) \cap C' &= \text{regiones } 2, 6 \text{ y } 7,
 \end{aligned}$$

y así sucesivamente.

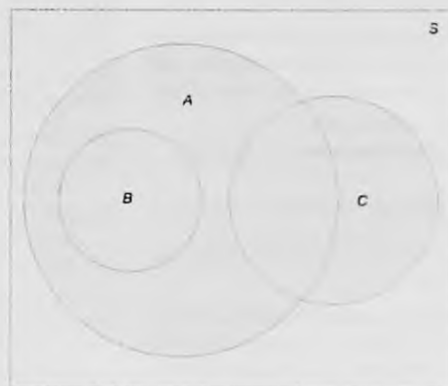


Figura 2.4: Eventos del espacio muestral S .

En la figura 2.4 vemos que los eventos A , B y C son subconjuntos del espacio muestral S . También es claro que el evento B es un subconjunto del evento A ; el evento $B \cap C$ no tiene elementos, por lo tanto, B y C son mutuamente excluyentes; el evento $A \cap C$ tiene al menos un elemento; y el evento $A \cup B = A$. Por consiguiente, la figura 2.4 podría representar una situación en la que se selecciona una carta al azar de una baraja ordinaria de 52 cartas y se observa si ocurren los siguientes eventos:

A : la carta es roja,

B : la carta es la jota, la reina o el rey de diamantes,

C : la carta es un as.

Claramente, el evento $A \cap C$ consta sólo de los dos ases rojos.

Varios resultados que se derivan de las definiciones precedentes, y que se pueden verificar de forma sencilla empleando diagramas de Venn, son como los que siguen:

- | | |
|---------------------------|---------------------------------|
| 1. $A \cap \phi = \phi$. | 6. $\phi' = S$. |
| 2. $A \cup \phi = A$. | 7. $(A')' = A$. |
| 3. $A \cap A' = \phi$. | 8. $(A \cap B)' = A' \cup B'$. |
| 4. $A \cup A' = S$. | 9. $(A \cup B)' = A' \cap B'$. |
| 5. $S' = \phi$. | |

Ejercicios

2.1 Liste los elementos de cada uno de los siguientes espacios muestrales:

- el conjunto de números enteros entre 1 y 50 que son divisibles entre 8;
- el conjunto $S = \{x \mid x^2 + 4x - 5 = 0\}$;
- el conjunto de resultados cuando se lanza una moneda al aire hasta que aparecen una cruz o tres caras;
- el conjunto $S = \{x \mid x \text{ es un continente}\}$;
- el conjunto $S = \{x \mid 2x - 4 \geq 0 \text{ y } x < 1\}$.

2.2 Utilice el método de la regla para describir el espacio muestral S , que consta de todos los puntos del primer cuadrante dentro de un círculo de radio 3 con centro en el origen.

2.3 ¿Cuáles de los siguientes eventos son iguales?

- $A = \{1, 3\}$;
- $B = \{x \mid x \text{ es un número de un dado}\}$;
- $C = \{x \mid x^2 - 4x + 3 = 0\}$;
- $D = \{x \mid x \text{ es el número de caras cuando se lanzan seis monedas al aire}\}$.

2.4 Un experimento implica lanzar un par de dados, uno verde y uno rojo, y registrar los números que resultan. Si x es igual al resultado en el dado verde y y es el resultado en el dado rojo, describa el espacio muestral S

- mediante la lista de los elementos (x, y) ;
- por medio del método de la regla.

2.5 Un experimento consiste en lanzar un dado y después lanzar una moneda una vez si el número en el dado es par. Si el número en el dado es impar, la moneda se lanza dos veces. Use la notación $4H$, por ejemplo, para denotar el resultado de que el dado muestre 4 y después la moneda caiga en cara, y $3HT$ para denotar el resultado de que el dado muestre 3, seguido por una cara y después una cruz en la moneda; construya un

diagrama de árbol para mostrar los 18 elementos del espacio muestral S .

2.6 De un grupo de cuatro suplentes se seleccionan dos jurados para servir en un juicio por homicidio. Utilice la notación A_1, A_3 , por ejemplo, para denotar el evento simple de que se seleccionen los suplentes 1 y 3, liste los 6 elementos del espacio muestral S .

2.7 De un grupo de estudiantes de química se seleccionan cuatro al azar y se clasifican como hombre o mujer. Liste los elementos del espacio muestral S_1 usando la letra H para hombre y M para mujer. Defina un segundo espacio muestral S_2 donde los elementos representen el número de mujeres seleccionadas.

2.8 Para el espacio muestral del ejercicio 2.4,

- liste los elementos que corresponden al evento A de que la suma sea mayor que 8;
- liste los elementos que corresponden al evento B de que ocurra un 2 en cualquiera de los dos dados;
- liste los elementos que corresponden al evento C de que salga un número mayor que 4 en el dado verde;
- liste los elementos que corresponden al evento $A \cap C$;
- liste los elementos que corresponden al evento $A \cap B$;
- liste los elementos que corresponden al evento $B \cap C$;
- construya un diagrama de Venn para ilustrar las intersecciones y uniones de los eventos A , B y C .

2.9 Para el espacio muestral del ejercicio 2.5,

- liste los elementos que corresponden al evento A en el que el dado salga un número menor que 3;
- liste los elementos que corresponden al evento B de que resulten 2 cruces;
- liste los elementos que corresponden al evento A' ;

- d) liste los elementos que corresponden al evento $A' \cap B$;
- e) liste los elementos que corresponden al evento $A \cup B$.

2.10 Se contrata a una empresa de ingenieros para que determine si ciertas vfas fluviales en Virginia, Estados Unidos, son seguras para la pesca. Se toman muestras de tres ríos.

- a) Liste los elementos de un espacio muestral S y utilice las letras P para "seguro para la pesca" y N para "inseguro para la pesca".
- b) Liste los elementos de S que correspondan al evento E de que al menos dos de los ríos son seguros para la pesca.
- c) Defina un evento que tiene como elementos a los puntos

{PPP, NPP, PPN, NPN}

2.11 El currículum de dos aspirantes masculinos para el puesto de profesor de química en una facultad se coloca en el mismo archivo que el de dos aspirantes mujeres. Hay dos puestos disponibles y el primero, con el rango de profesor asistente, se cubre seleccionando al azar a uno de los cuatro aspirantes. El segundo puesto, con el rango de profesor titular, se cubre después mediante la selección aleatoria de uno de los tres aspirantes restantes. Utilice la notación H_2M_1 , por ejemplo, para denotar el evento simple de que el primer puesto se cubra con el segundo aspirante hombre y el segundo puesto se cubra después con la primera aspirante mujer,

- a) liste los elementos de un espacio muestral S ;
- b) liste los elementos de S que corresponden al evento A en que el puesto de profesor asistente se cubre con un aspirante hombre;
- c) liste los elementos de S que corresponden al evento B en que exactamente 1 de los 2 puestos se cubre con un aspirante hombre;
- d) liste los elementos de S que corresponden al evento C en que ningún puesto se cubre con un aspirante hombre;
- e) liste los elementos de S que corresponden al evento $A \cap B$;
- f) liste los elementos de S que corresponden al evento $A \cup C$;
- g) construya un diagrama de Venn para ilustrar las intersecciones y las uniones de los eventos A , B y C .

2.12 Se estudian el ejercicio y la dieta como posibles sustitutos del medicamento para bajar la presión sanguínea. Se utilizarán tres grupos de individuos para estudiar el efecto del ejercicio. Los integrantes del grupo uno son sedentarios, los del dos caminan y los del tres nadan una hora al día. La mitad de cada uno de los tres grupos de ejercicio tendrá una dieta sin sal. Un gru-

po adicional de individuos no hará ejercicio ni restringirá su consumo de sal, pero tomará el medicamento estándar. Use Z para sedentario, C para caminante, S para nadador, Y para sal, N para sin sal, M para medicamento y F para sin medicamento.

- a) Muestre todos los elementos del espacio muestral S .
- b) Dado que A es el conjunto de individuos sin medicamento y B es el conjunto de caminantes, liste los elementos de $A \cup B$.
- c) Liste los elementos de $A \cap B$.

2.13 Construya un diagrama de Venn para ilustrar las posibles intersecciones y uniones en los siguientes eventos relativos al espacio muestral que consta de todos los automóviles fabricados en Estados Unidos.

C : cuatro puertas, T : techo corredizo, D : dirección hidráulica

2.14 Si $S = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ y $A = \{0, 2, 4, 6, 8\}$, $B = \{1, 3, 5, 7, 9\}$, $C = \{2, 3, 4, 5\}$ y $D = \{1, 6, 7\}$, liste los elementos de los conjuntos que corresponden a los siguientes eventos:

- a) $A \cup C$;
- b) $A \cap B$;
- c) C' ;
- d) $(C' \cap D) \cup B$;
- e) $(S \cap C)'$;
- f) $A \cap C \cap D'$.

2.15 Considere el espacio muestral $S = \{\text{cobre, sodio, nitrógeno, potasio, uranio, oxígeno, cinc}\}$ y los eventos

$$A = \{\text{cobre, sodio, cinc}\},$$

$$B = \{\text{sodio, nitrógeno, potasio}\}$$

$$C = \{\text{oxígeno}\}.$$

Liste los elementos de los conjuntos que corresponden a los siguientes eventos:

- a) A' ;
- b) $A \cup C$;
- c) $(A \cap B') \cup C'$;
- d) $B' \cap C'$;
- e) $A \cap B \cap C$;
- f) $(A' \cup B') \cap (A \cap C)$.

2.16 Si $S = \{x \mid 0 < x < 12\}$, $M = \{x \mid 1 < x < 9\}$ y $N = \{x \mid 0 < x < 5\}$, encuentre

- a) $M \cup N$;
- b) $M \cap N$;
- c) $M' \cap N'$.

2.17 Sean A , B y C eventos relativos al espacio muestral S . Utilice diagramas de Venn para sombrear las áreas que representan los siguientes eventos:

- a) $(A \cap B)'$;
- b) $(A \cup B)'$;
- c) $(A \cap C) \cup B$.

2.18 ¿Cuál de los siguientes pares de eventos son mutuamente excluyentes?

- Un golfista que se clasifica en último lugar en la vuelta del hoyo 18, en un torneo de 72 hoyos, y pierde el torneo.
- Un jugador de póquer que tiene flor (todas las cartas del mismo palo) y 3 del mismo palo en la misma mano de 5 cartas.
- Una madre que da a luz a una niña y a un par de gemelas el mismo día.
- Un jugador de ajedrez que pierde el último juego y gana el torneo.

2.19 Suponga que una familia sale de vacaciones de verano en su casa rodante y que M es el evento de que sufrirán fallas mecánicas, T es el evento de que recibirán una infracción por cometer una falta de tránsito y V es el evento de que llegarán a un lugar para acampar que esté lleno. Remítase al diagrama de Venn de la figura 2.5 y exprese con palabras los eventos representados por las siguientes regiones:

- región 5;
- región 3;
- regiones 1 y 2 juntas;
- regiones 4 y 7 juntas;
- regiones 3, 6, 7 y 8 juntas.

2.20 Remítase al ejercicio 2.19 y al diagrama de Venn de la figura 2.5, liste los números de las regiones que representan los siguientes eventos:

- La familia no experimentará fallas mecánicas y no será multada por cometer una infracción de tránsito, pero llegará a un lugar para acampar que está lleno.
- La familia experimentará tanto fallas mecánicas como problemas para localizar un lugar disponible para acampar, pero no será multada por cometer una infracción de tránsito.
- La familia experimentará fallas mecánicas o encontrará un lugar para acampar lleno, pero no será multada por cometer una infracción de tránsito.
- La familia no llegará a un lugar para acampar lleno.



Figura 2.5: Diagrama de Venn para los ejercicios 2.19 y 2.20.

2.3 Conteo de puntos muestrales

Uno de los problemas que el estadístico debe considerar e intentar evaluar es el elemento de aleatoriedad asociado con la ocurrencia de ciertos eventos cuando se realiza un experimento. Estos problemas pertenecen al campo de la probabilidad, un tema que se estudiará en la sección 2.4. En muchos casos debemos ser capaces de resolver un problema de probabilidad mediante el conteo del número de puntos en el espacio muestral, sin listar realmente cada elemento. El principio fundamental del conteo, a menudo denominado **regla de multiplicación**, se establece en la regla 2.1.

Regla 2.1: Si una operación se puede llevar a cabo en n_1 formas, y si para cada una de éstas se puede realizar una segunda operación en n_2 formas, entonces las dos operaciones se pueden ejecutar juntas de $n_1 n_2$ formas.

Ejemplo 2.13: ¿Cuántos puntos muestrales hay en el espacio muestral cuando se lanza un par de dados una vez?

Solución: El primer dado puede caer en cualquiera de $n_1 = 6$ maneras. Para cada una de esas 6 maneras el segundo dado también puede caer en $n_2 = 6$ formas. Por lo tanto, el par de dados puede caer en $n_1 n_2 = (6)(6) = 36$ formas posibles. ▮

Ejemplo 2.14: Un urbanista de una nueva subdivisión ofrece a los posibles compradores de una casa elegir entre Tudor, rústica, colonial y tradicional el estilo de la fachada, y entre una planta, dos pisos y desniveles el plano de construcción. ¿En cuántas formas diferentes puede un comprador ordenar una de estas casas?

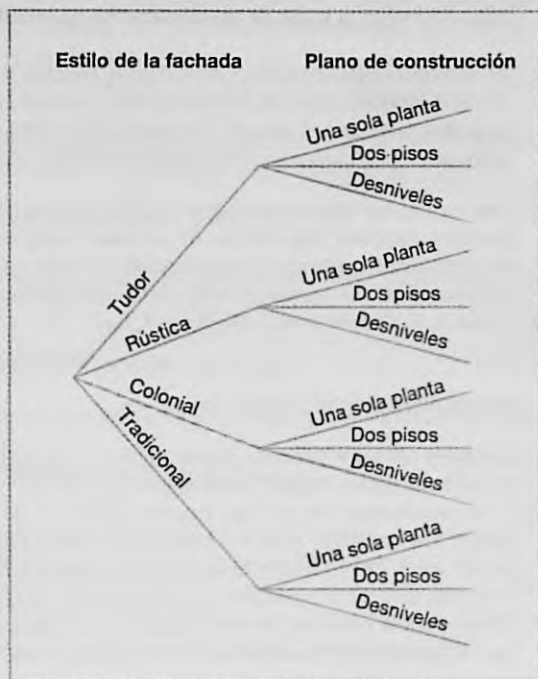


Figura 2.6: Diagrama de árbol para el ejemplo 2.14.

Solución: Como $n_1 = 4$ y $n_2 = 3$, un comprador debe elegir entre

$$n_1 n_2 = (4)(3) = 12 \text{ casas posibles.} \quad \blacksquare$$

Las respuestas a los dos ejemplos anteriores se comprueban construyendo diagramas de árbol y contando las diversas trayectorias a lo largo de las ramas. Así, en el

ejemplo 2.14 habrá $n_1 = 4$ ramas que corresponden a los diferentes estilos de la fachada, y después habrá $n_2 = 3$ ramas que se extienden de cada una de estas 4 ramas para representar los diferentes planos de plantas. Este diagrama de árbol, como se ilustra en la figura 2.6, proporciona las $n_1 n_2 = 12$ opciones de casas dadas por las trayectorias a lo largo de las ramas.

Ejemplo 2.15: Si un miembro de un club que tiene 22 integrantes necesitara elegir un presidente y un tesorero, ¿de cuántas maneras diferentes se podría elegir a ambos?

Solución: Para el puesto de presidente hay 22 posibilidades en total. Para cada una de esas 22 posibilidades hay 21 posibilidades de elegir al tesorero. Si utilizamos la regla de la multiplicación, obtenemos $n_1 \times n_2 = 22 \times 21 = 462$ maneras diferentes. ─

La regla de la multiplicación (regla 2.1) se puede extender para abarcar cualquier número de operaciones. Por ejemplo, suponga que un cliente desea comprar un nuevo teléfono celular y que puede elegir entre $n_1 = 5$ marcas, $n_2 = 5$ tipos de capacidad y $n_3 = 4$ colores. Estas tres clasificaciones dan como resultado $n_1 n_2 n_3 = (5)(5)(4) = 100$ diferentes formas en las que un cliente puede ordenar uno de estos teléfonos. A continuación se formula la **regla de multiplicación generalizada** que cubre k operaciones.

Regla 2.2: Si una operación se puede ejecutar en n_1 formas, y si para cada una de éstas se puede llevar a cabo una segunda operación en n_2 formas, y para cada una de las primeras dos se puede realizar una tercera operación en n_3 formas, y así sucesivamente, entonces la serie de k operaciones se puede realizar en $n_1 n_2 \dots n_k$ formas.

Ejemplo 2.16: Sam va a armar una computadora y para comprar las partes tiene que elegir entre las siguientes opciones: dos marcas de circuitos integrados, cuatro marcas de discos duros, tres marcas de memorias y cinco tiendas locales en las que puede adquirir un conjunto de accesorios. ¿De cuántas formas diferentes puede Sam comprar las partes?

Solución: Como $n_1 = 2$, $n_2 = 4$, $n_3 = 3$ y $n_4 = 5$, hay

$$n_1 \times n_2 \times n_3 \times n_4 = 2 \times 4 \times 3 \times 5 = 120$$

formas diferentes de comprar las partes. ─

Ejemplo 2.17: ¿Cuántos números pares de cuatro dígitos se pueden formar con los dígitos 0, 1, 2, 5, 6 y 9, si cada dígito se puede usar sólo una vez?

Solución: Como el número debe ser par, tenemos sólo $n_1 = 3$ opciones para la posición de las unidades. Sin embargo, para un número de cuatro dígitos la posición de los millares no puede ser 0. Por lo tanto, consideramos la posición de las unidades en dos partes: 0 o diferente de 0. Si la posición de las unidades es 0 (es decir, $n_1 = 1$), tenemos $n_2 = 5$ opciones para la posición de los millares, $n_3 = 4$ para la posición de las centenas y $n_4 = 3$ para la posición de las decenas. Por lo tanto, en este caso tenemos un total de

$$n_1 n_2 n_3 n_4 = (1)(5)(4)(3) = 60$$

números pares de cuatro dígitos. Por otro lado, si la posición de las unidades no es 0 (es decir, $n_1 = 2$), tenemos $n_2 = 4$ opciones para la posición de los millares, $n_3 = 4$ para la posición de las centenas y $n_4 = 3$ para la posición de las decenas. En esta situación tenemos un total de

$$n_1 n_2 n_3 n_4 = (2)(4)(4)(3) = 96$$

números pares de cuatro dígitos.

Puesto que los dos casos anteriores son mutuamente excluyentes, el número total de números pares de cuatro dígitos se puede calcular usando $60 + 96 = 156$. ▮

Con frecuencia nos interesamos en un espacio muestral que contiene como elementos a todas las posibles ordenaciones o arreglos de un grupo de objetos. Por ejemplo, cuando queremos saber cuántos arreglos diferentes son posibles para sentar a seis personas alrededor de una mesa, o cuando nos preguntamos cuántas ordenaciones diferentes son posibles para sacar dos billetes de lotería de un total de 20. En este caso los diferentes arreglos se llaman **permutaciones**.

Definición 2.7: Una **permutación** es un arreglo de todo o parte de un conjunto de objetos.

Considere las tres letras a, b y c . Las permutaciones posibles son abc, acb, bac, bca, cab y cba , por lo tanto, vemos que hay 6 arreglos distintos. Si utilizamos la regla 2.2 podemos llegar a la respuesta 6 sin listar realmente las diferentes ordenaciones. Hay $n_1 = 3$ opciones para la primera posición. Sin importar cuál letra se elija, siempre habrá $n_2 = 2$ opciones para la segunda posición. Por último, independientemente de cuál de las dos letras se elija para las primeras dos posiciones, sólo hay $n_3 = 1$ elección para la última posición, lo que da un total de

$$n_1 n_2 n_3 = (3)(2)(1) = 6 \text{ permutaciones}$$

mediante la regla 2.2. En general, n objetos distintos se pueden arreglar en

$$n(n-1)(n-2) \dots (3)(2)(1) \text{ formas.}$$

Existe una notación para una cifra como ésta.

Definición 2.8 Para cualquier entero no negativo n , $n!$, denominado “ n factorial” se define como

$$N! = n(n-1) \dots (2)(1),$$

con el caso especial de $0! = 1$.

Si utilizamos el argumento anterior llegamos al siguiente teorema.

Teorema 2.1: El número de permutaciones de n objetos es $n!$

El número de permutaciones de las cuatro letras a, b, c y d será $4! = 24$. Consideremos ahora el número de permutaciones que son posibles tomando dos de las cuatro letras a la vez. Éstas serían $ab, ac, ad, ba, bc, bd, ca, cb, cd, da, db$ y dc . De nuevo, si utilizamos la regla 2.1, tenemos dos posiciones para llenar con $n_1 = 4$ opciones para la primera y después $n_2 = 3$ opciones para la segunda, para un total de

$$n_1 n_2 = (4)(3) = 12$$

permutaciones. En general, n objetos distintos tomados de r a la vez se pueden arreglar en

$$n(n-1)(n-2) \dots (n-r+1)$$

formas. Representamos este producto mediante

$${}_n P_r = \frac{n!}{(n-r)!}.$$

Como resultado tenemos el teorema que sigue.

Teorema 2.2: El número de permutaciones de n objetos distintos tomados de r a la vez es

$${}_n P_r = \frac{n!}{(n-r)!}$$

Ejemplo 2.18: En un año se otorgará uno de tres premios (a la investigación, la enseñanza y el servicio) a algunos de los estudiantes, de un grupo de 25, de posgrado del departamento de estadística. Si cada estudiante puede recibir un premio como máximo, ¿cuántas selecciones posibles habría?

Solución: Como los premios son distinguibles, se trata de un problema de permutación. El número total de puntos muestrales es

$${}_{25} P_3 = \frac{25!}{(25-3)!} = \frac{25!}{22!} = (25)(24)(23) = 13,800.$$

Ejemplo 2.19: En un club estudiantil compuesto por 50 personas se va a elegir a un presidente y a un tesorero. ¿Cuántas opciones diferentes de funcionarios son posibles si

- no hay restricciones;
- A participará sólo si él es el presidente;
- B y C participarán juntos o no lo harán;
- D y E no participarán juntos?

Solución: a) El número total de opciones de funcionarios, si no hay restricciones, es

$${}_{50} P_2 = \frac{50!}{48!} = (50)(49) = 2450.$$

- Como A participaría sólo si es el presidente, tenemos dos situaciones: i) A se elige como presidente, lo cual produce 49 resultados posibles para el puesto de tesorero; o ii) los funcionarios se eligen de entre las 49 personas restantes sin tomar en cuenta a A, en cuyo caso el número de opciones es ${}_{49} P_2 = (49)(48) = 2352$. Por lo tanto, el número total de opciones es $49 + 2352 = 2401$.
- El número de selecciones cuando B y C participan juntos es 2. El número de selecciones cuando ni B ni C se eligen es ${}_{48} P_2 = 2256$. Por lo tanto, el número total de opciones en esta situación es $2 + 2256 = 2258$.
- El número de selecciones cuando D participa como funcionario pero sin E es $(2)(48) = 96$, donde 2 es el número de puestos que D puede ocupar y 48 es el número de selecciones de los otros funcionarios de las personas restantes en el club, excepto E. El número de selecciones cuando E participa como funcionario pero sin D también es $(2)(48) = 96$. El número de selecciones cuando tanto D como E no son elegidos es ${}_{48} P_2 = 2256$. Por lo tanto, el número total de opciones es $(2)(96) + 2256 = 2448$. Este problema también tiene otra solución rápida: como D y E sólo pueden participar juntos de dos maneras, la respuesta es $2450 - 2 = 2448$. ■

Las permutaciones que ocurren al arreglar objetos en un círculo se llaman **permutaciones circulares**. Dos permutaciones circulares no se consideran diferentes a menos que los objetos correspondientes en los dos arreglos estén precedidos o seguidos por un objeto diferente, conforme avancemos en la dirección de las manecillas del reloj. Por ejemplo, si cuatro personas juegan *bridge*, no tenemos una permutación nueva si se mueven una posición en la dirección de las manecillas del reloj. Si consideramos a una persona en una posición fija y arreglamos a las otras tres de $3!$ formas, encontramos que hay seis arreglos distintos para el juego de *bridge*.

Teorema 2.3: El número de permutaciones de n objetos ordenados en un círculo es $(n - 1)!$.

Hasta ahora hemos considerado permutaciones de objetos distintos. Es decir, todos los objetos fueron por completo diferentes o distinguibles. Evidentemente, si tanto la letra b como la c son iguales a x , entonces las 6 permutaciones de las letras a , b y c se convierten en axx , axx , xax , xax , xxa y xxa , de las cuales sólo 3 son diferentes. Por lo tanto, con 3 letras, en las que 2 son iguales, tenemos $3!/2! = 3$ permutaciones distintas. Con 4 letras diferentes a , b , c y d tenemos 24 permutaciones distintas. Si permitimos que $a = b = x$ y $c = d = y$, podemos listar sólo las siguientes permutaciones distintas: $xyxy$, $xyxy$, $yxyx$, $yxyx$, $xyyx$ y $yxyx$. De esta forma tenemos $4!/(2!2!) = 6$ permutaciones distintas.

Teorema 2.4: El número de permutaciones distintas de n objetos, en el que n_1 son de una clase, n_2 de una segunda clase, ..., n_k de una k -ésima clase es

$$\frac{n!}{n_1! n_2! \cdots n_k!}$$

Ejemplo 2.20: Durante un entrenamiento de fútbol americano colegial, el coordinador defensivo necesita tener a 10 jugadores parados en una fila. Entre estos 10 jugadores hay 1 de primer año, 2 de segundo año, 4 de tercer año y 3 de cuarto año, respectivamente. ¿De cuántas formas diferentes se pueden arreglar en una fila si lo único que los distingue es el grado en el cual están?

Solución: Usando directamente el teorema 2.4, el número total de arreglos es

$$\frac{10!}{1! 2! 4! 3!} = 12,600.$$

Con frecuencia nos interesa el número de formas de dividir un conjunto de n objetos en r subconjuntos denominados **celdas**. Se consigue una partición si la intersección de todo par posible de los r subconjuntos es el conjunto vacío ϕ , y si la unión de todos los subconjuntos da el conjunto original. El orden de los elementos dentro de una celda no tiene importancia. Considere el conjunto $\{a, e, i, o, u\}$. Las particiones posibles en dos celdas en las que la primera celda contenga 4 elementos y la segunda 1 son

$$\{(a, e, i, o), (u)\}, \{(a, i, o, u), (e)\}, \{(e, i, o, u), (a)\}, \{(a, e, o, u), (i)\}, \{(a, e, i, u), (o)\}.$$

Vemos que hay 5 formas de partir un conjunto de 4 elementos en dos subconjuntos o celdas que contengan 4 elementos en la primera celda y 1 en la segunda.

El número de particiones para esta ilustración se denota con la expresión

$$\binom{5}{4, 1} = \frac{5!}{4! 1!} = 5,$$

en la que el número superior representa el número total de elementos y los números inferiores representan el número de elementos que van en cada celda. Establecemos esto de forma más general en el teorema 2.5.

Teorema 2.5: El número de formas de partir un conjunto de n objetos en r celdas con n_1 elementos en la primera celda, n_2 elementos en la segunda, y así sucesivamente, es

$$\binom{n}{n_1, n_2, \dots, n_r} = \frac{n!}{n_1! n_2! \cdots n_r!},$$

donde $n_1 + n_2 + \dots + n_r = n$.

Ejemplo 2.21: Un hotel va a hospedar a siete estudiantes de posgrado que asisten a una conferencia, ¿en cuántas formas los puede asignar a una habitación triple y a dos dobles?

Solución: El número total de particiones posibles sería

$$\binom{7}{3, 2, 2} = \frac{7!}{3! 2! 2!} = 210. \quad \blacksquare$$

En muchos problemas nos interesamos en el número de formas de seleccionar r objetos de n sin importar el orden. Tales selecciones se llaman **combinaciones**. Una combinación es realmente una partición con dos celdas, donde una celda contiene los r objetos seleccionados y la otra contiene los $(n - r)$ objetos restantes. El número de tales combinaciones se denota con

$$\binom{n}{r, n - r}, \text{ que por lo general se reduce a } \binom{n}{r},$$

debido a que el número de elementos en la segunda celda debe ser $n - r$.

Teorema 2.6: El número de combinaciones de n objetos distintos tomados de r a la vez es

$$\binom{n}{r} = \frac{n!}{r!(n - r)!}.$$

Ejemplo 2.22: Un niño le pide a su madre que le lleve cinco cartuchos de Game-Boy™ de su colección de 10 juegos recreativos y 5 de deportes. ¿De cuántas maneras podría su madre llevarle 3 juegos recreativos y 2 de deportes?

Solución: El número de formas de seleccionar 3 cartuchos de 10 es

$$\binom{10}{3} = \frac{10!}{3!(10 - 3)!} = 120.$$

El número de formas de seleccionar 2 cartuchos de 5 es

$$\binom{5}{2} = \frac{5!}{2! 3!} = 10.$$

Si utilizamos la regla de la multiplicación (regla 2.1) con $n_1 = 120$ y $n_2 = 10$, tenemos que hay $(120)(10) = 1200$ formas. ■

Ejemplo 2.23: ¿Cuántos arreglos diferentes de letras se pueden hacer con las letras de la palabra *STATISTICS*?

Solución: Si utilizamos el mismo argumento expuesto en el teorema 2.6, en este ejemplo podemos realmente aplicar el teorema 2.5 para obtener

$$\binom{10}{3, 3, 2, 1, 1} = \frac{10!}{3! 3! 2! 1! 1!} = 50,400.$$

Aquí tenemos 10 letras en total, donde 2 letras (*S*, *T*) aparecen tres veces cada una, la letra *I* aparece dos veces, y las letras *A* y *C* aparecen una vez cada una. Por otro lado, el resultado se puede obtener directamente usando el teorema 2.4. ■

Ejercicios

2.21 A los participantes de una convención se les ofrecen seis recorridos, cada uno de tres días, a sitios de interés. ¿De cuántas maneras se puede acomodar una persona para que vaya a uno de los recorridos planeados por la convención?

2.22 En un estudio médico los pacientes se clasifican en 8 formas de acuerdo con su tipo sanguíneo: AB^+ , AB^- , A^+ , A^- , B^+ , B^- , O^+ u O^- ; y también de acuerdo con su presión sanguínea: baja, normal o alta. Encuentre el número de formas en las que se puede clasificar a un paciente.

2.23 Si un experimento consiste en lanzar un dado y después extraer una letra al azar del alfabeto inglés, ¿cuántos puntos habrá en el espacio muestral?

2.24 Los estudiantes de humanidades de una universidad privada se clasifican como estudiantes de primer año, de segundo año, de penúltimo año o de último año, y también de acuerdo con su género (hombres o mujeres). Calcule el número total de clasificaciones posibles para los estudiantes de esa universidad.

2.25 Cierta marca de calzado existe en 5 diferentes estilos y cada estilo está disponible en 4 colores distintos. Si la tienda deseara mostrar la cantidad de pares de zapatos que incluya todos los diversos estilos y colores, ¿cuántos pares diferentes tendría que mostrar?

2.26 Un estudio en California concluyó que siguiendo siete sencillas reglas para la salud un hombre y una mujer pueden prolongar su vida 11 y 7 años en promedio, respectivamente. Estas 7 reglas son: no fumar, hacer ejercicio de manera habitual, moderar su consumo de alcohol, dormir siete u ocho horas, mantener el peso adecuado, desayunar y no ingerir alimentos entre comi-

das. De cuántas formas puede una persona adoptar cinco de estas reglas:

- a) ¿Si la persona actualmente infringe las siete reglas?
b) ¿Si la persona nunca bebe y siempre desayuna?

2.27 Un urbanista de un nuevo fraccionamiento ofrece a un posible comprador de una casa elegir entre 4 diseños, 3 diferentes sistemas de calefacción, un garaje o cobertizo, y un patio o un porche cubierto. ¿De cuántos planos diferentes dispone el comprador?

2.28 Un medicamento para aliviar el asma se puede adquirir en 5 diferentes laboratorios y en forma de líquido, comprimidos o cápsulas, todas en concentración normal o alta. ¿De cuántas formas diferentes puede un médico recetar la medicina a un paciente que sufre de asma?

2.29 En un estudio económico de combustibles, cada uno de 3 autos de carreras se prueba con 5 marcas diferentes de gasolina en 7 lugares de prueba que se localizan en diferentes regiones del país. Si en el estudio se utilizan 2 pilotos y las pruebas se realizan una vez en cada uno de los distintos grupos de condiciones, ¿cuántas pruebas se necesita realizar?

2.30 ¿De cuántas formas distintas se puede responder una prueba de falso-verdadero que consta de 9 preguntas?

2.31 Un testigo de un accidente automovilístico le dijo a la policía que la matrícula del culpable, que huyó, contenía las letras RLH seguidas por 3 dígitos, de los cuales el primero era un 5. Si el testigo no recuerda los 2 últimos dígitos, pero está seguro de que los 3 eran distintos, calcule la cantidad máxima de registros de automóviles que la policía tendría que revisar.

- 2.32** a) ¿De cuántas maneras se pueden formar 6 personas para abordar un autobús?
 b) ¿Cuántas maneras son posibles si, de las 6, 3 personas específicas insisten en formarse una después de la otra?
 c) ¿De cuántas maneras se pueden formar si, de las 6, 2 personas específicas se rehúsan a formarse una detrás de la otra?
- 2.33** Si una prueba de opción múltiple consta de 5 preguntas, cada una con 4 respuestas posibles, de las cuales sólo 1 es correcta,
 a) ¿de cuántas formas diferentes puede un estudiante elegir una respuesta a cada pregunta?
 b) ¿de cuántas maneras puede un estudiante elegir una respuesta a cada pregunta y obtener todas las respuestas incorrectas?
- 2.34** a) ¿Cuántas permutaciones distintas se pueden hacer con las letras de la palabra *COLUMNA*?
 b) ¿Cuántas de estas permutaciones comienzan con la letra *M*?
- 2.35** Un contratista desea construir 9 casas, cada una con diferente diseño. ¿De cuántas formas puede ubicarlas en la calle en la que las va a construir si en un lado de ésta hay 6 lotes y en el lado opuesto hay 3?
- 2.36** a) ¿Cuántos números de tres dígitos se pueden formar con los dígitos 0, 1, 2, 3, 4, 5 y 6 si cada dígito se puede usar sólo una vez?
 b) ¿Cuántos de estos números son impares?
 c) ¿Cuántos son mayores que 330?
- 2.37** ¿De cuántas maneras se pueden sentar 4 niños y 5 niñas en una fila, si se deben alternar unos y otras?
- 2.38** Cuatro parejas compran 8 lugares en la misma fila para un concierto. ¿De cuántas maneras diferentes se pueden sentar...
 a) sin restricciones?
 b) si cada pareja se sienta junta?
 c) si todos los hombres se sientan juntos a la derecha de todas las mujeres?
- 2.39** En un concurso regional de ortografía, los 8 finalistas son 3 niños y 5 niñas. Encuentre el número de puntos muestrales en el espacio muestral S para el número de ordenamientos posibles al final del concurso para
 a) los 8 finalistas;
 b) los 3 primeros lugares.
- 2.40** ¿De cuántas formas se pueden cubrir las 5 posiciones iniciales en un equipo de baloncesto con 8 jugadores que pueden jugar cualquiera de las posiciones?
- 2.41** Encuentre el número de formas en que se puede asignar 6 profesores a 4 secciones de un curso introductorio de psicología, si ningún profesor se asigna a más de una sección.
- 2.42** De un grupo de 40 boletos se sacan 3 billetes de lotería para el primero, segundo y tercer premios. Encuentre el número de puntos muestrales en S para dar los 3 premios, si cada concursante sólo tiene un boleto.
- 2.43** ¿De cuántas maneras se pueden plantar 5 árboles diferentes en un círculo?
- 2.44** ¿De cuántas formas se puede acomodar en círculo una caravana de ocho carretas de Arizona?
- 2.45** ¿Cuántas permutaciones distintas se pueden hacer con las letras de la palabra *INFINITO*?
- 2.46** ¿De cuántas maneras se pueden colocar 3 robles, 4 pinos y 2 arces a lo largo de la línea divisoria de una propiedad, si no se distingue entre árboles del mismo tipo?
- 2.47** ¿De cuántas formas se puede seleccionar a 3 de 8 candidatos recién graduados, igualmente calificados, para ocupar las vacantes de un despacho de contabilidad?
- 2.48** ¿Cuántas formas hay en que dos estudiantes no tengan la misma fecha de cumpleaños en un grupo de 60?

2.4 Probabilidad de un evento

Quizá fue la insaciable sed del ser humano por el juego lo que condujo al desarrollo temprano de la teoría de la probabilidad. En un esfuerzo por aumentar sus triunfos, algunos pidieron a los matemáticos que les proporcionaran las estrategias óptimas para los diversos juegos de azar. Algunos de los matemáticos que brindaron tales estrategias fueron Pascal, Leibniz, Fermat y James Bernoulli. Como resultado de este desarrollo inicial de la teoría de la probabilidad, la inferencia estadística, con todas sus predicciones y generalizaciones, ha rebasado el ámbito de los juegos de azar para abarcar muchos otros campos asociados con los eventos aleatorios, como la política, los negocios, el pronóstico del clima y la

investigación científica. Para que estas predicciones y generalizaciones sean razonablemente precisas, resulta esencial la comprensión de la teoría básica de la probabilidad.

¿A qué nos referimos cuando hacemos afirmaciones como “Juan probablemente ganará el torneo de tenis”, o “tengo 50% de probabilidades de obtener un número par cuando lanzo un dado”, o “la universidad no tiene posibilidades de ganar el juego de fútbol esta noche”, o “la mayoría de nuestros graduados probablemente estarán casados dentro de tres años”? En cada caso expresamos un resultado del cual no estamos seguros, pero con base en la experiencia, o a partir de la comprensión de la estructura del experimento, confiamos hasta cierto punto en la validez de nuestra afirmación.

En el resto de este capítulo consideraremos sólo aquellos experimentos para los cuales el espacio muestral contiene un número finito de elementos. La probabilidad de la ocurrencia de un evento que resulta de tal experimento estadístico se evalúa utilizando un conjunto de números reales denominados **pesos** o **probabilidades**, que van de 0 a 1. Para todo punto en el espacio muestral asignamos una probabilidad tal que la suma de todas las probabilidades es 1. Si tenemos razón para creer que al llevar a cabo el experimento es bastante probable que ocurra cierto punto muestral, le tendríamos que asignar a éste una probabilidad cercana a 1. Por el contrario, si creemos que no hay probabilidades de que ocurra cierto punto muestral, le tendríamos que asignar a éste una probabilidad cercana a cero. En muchos experimentos, como lanzar una moneda o un dado, todos los puntos muestrales tienen la misma oportunidad de ocurrencia, por lo tanto, se les asignan probabilidades iguales. A los puntos fuera del espacio muestral, es decir, a los eventos simples que no tienen posibilidades de ocurrir, les asignamos una probabilidad de cero.

Para encontrar la probabilidad de un evento A sumamos todas las probabilidades que se asignan a los puntos muestrales en A . Esta suma se denomina **probabilidad** de A y se denota con $P(A)$.

Definición 2.9: La **probabilidad** de un evento A es la suma de los pesos de todos los puntos muestrales en A . Por lo tanto,

$$0 \leq P(A) \leq 1, \quad P(\phi) = 0 \quad \text{y} \quad P(S) = 1.$$

Además, si A_1, A_2, A_3, \dots es una serie de eventos mutuamente excluyentes, entonces

$$P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$$

Ejemplo 2.24 Una moneda se lanza dos veces. ¿Cuál es la probabilidad de que ocurra al menos una cara (H)?

Solución: El espacio muestral para este experimento es

$$S = \{HH, HT, TH, TT\}$$

Si la moneda está balanceada, cada uno de estos resultados tendrá las mismas probabilidades de ocurrir. Por lo tanto, asignamos una probabilidad de ω a cada uno de los puntos muestrales. Entonces, $4\omega = 1$ o $\omega = 1/4$. Si A representa el evento de que ocurra al menos una cara (H), entonces

$$A = \{HH, HT, TH\} \text{ y } P(A) = \frac{1}{4} + \frac{1}{4} + \frac{1}{4} = \frac{3}{4}.$$

Ejemplo 2.25: Se carga un dado de forma que exista el doble de probabilidades de que salga un número par que uno impar. Si E es el evento de que ocurra un número menor que 4 en un solo lanzamiento del dado, calcule $P(E)$.

Solución: El espacio muestral es $S = \{1, 2, 3, 4, 5, 6\}$. Asignamos una probabilidad de w a cada número impar y una probabilidad de $2w$ a cada número par. Como la suma de las probabilidades debe ser 1, tenemos $9w = 1$ o $w = 1/9$. Por lo tanto, asignamos probabilidades de $1/9$ y $2/9$ a cada número impar y par, respectivamente. Por consiguiente,

$$E = \{1, 2, 3\} \text{ y } P(E) = \frac{1}{9} + \frac{2}{9} + \frac{1}{9} = \frac{4}{9}.$$

Ejemplo 2.26: En el ejemplo 2.25, sea A el evento de que resulte un número par y sea B el evento de que resulte un número divisible entre 3. Calcule $P(A \cup B)$ y $P(A \cap B)$.

Solución: Para los eventos $A = \{2, 4, 6\}$ y $B = \{3, 6\}$, tenemos

$$A \cup B = \{2, 3, 4, 6\} \text{ y } A \cap B = \{6\}.$$

Al asignar una probabilidad de $1/9$ a cada número impar y de $2/9$ a cada número par, tenemos

$$P(A \cup B) = \frac{2}{9} + \frac{1}{9} + \frac{2}{9} + \frac{2}{9} = \frac{7}{9} \text{ y } P(A \cap B) = \frac{2}{9}.$$

Si el espacio muestral para un experimento contiene N elementos, todos los cuales tienen las mismas probabilidades de ocurrir, asignamos una probabilidad igual a $1/N$ a cada uno de los N puntos. La probabilidad de que cualquier evento A contenga n de estos N puntos muestrales es entonces el cociente del número de elementos en A y el número de elementos en S .

Regla 2.3: Si un experimento puede dar como resultado cualquiera de N diferentes resultados que tienen las mismas probabilidades de ocurrir, y si exactamente n de estos resultados corresponden al evento A , entonces la probabilidad del evento A es

$$P(A) = \frac{n}{N}.$$

Ejemplo 2.27: A una clase de estadística para ingenieros asisten 25 estudiantes de ingeniería industrial, 10 de ingeniería mecánica, 10 de ingeniería eléctrica y 8 de ingeniería civil. Si el profesor elige al azar a un estudiante para que conteste una pregunta, ¿qué probabilidades hay de que el elegido sea a) estudiante de ingeniería industrial, b) estudiante de ingeniería civil o estudiante de ingeniería eléctrica?

Solución: Las especialidades de los estudiantes de ingeniería industrial, mecánica, eléctrica y civil se denotan con I , M , E y C , respectivamente. El grupo está integrado por 53 estudiantes y todos tienen las mismas probabilidades de ser seleccionados.

a) Como 25 de los 53 individuos estudian ingeniería industrial, la probabilidad del evento I , es decir, la de elegir al azar a alguien que estudia ingeniería industrial, es

$$P(I) = \frac{25}{53}.$$

b) Como 18 de los 53 estudiantes son de las especialidades de ingeniería civil o eléctrica, se deduce que

$$P(C \cup E) = \frac{18}{53}.$$

Ejemplo 2.28: En una mano de póquer que consta de 5 cartas encuentre la probabilidad de tener 2 ases y 3 jotas.

Solución: El número de formas de tener 2 ases de 4 cartas es

$$\binom{4}{2} = \frac{4!}{2! 2!} = 6,$$

y el número de formas de tener 3 jotas de 4 cartas es

$$\binom{4}{3} = \frac{4!}{3! 1!} = 4.$$

Mediante la regla de multiplicación (regla 2.1), obtenemos $n = (6)(4) = 24$ manos con 2 ases y 3 jotas. El número total de manos de póquer de 5 cartas, todas las cuales tienen las mismas probabilidades de ocurrir, es

$$N = \binom{52}{5} = \frac{52!}{5! 47!} = 2,598,960.$$

Por lo tanto, la probabilidad del evento C de obtener 2 ases y 3 jotas en una mano de póquer de 5 cartas es

$$P(C) = \frac{24}{2,598,960} = 0.9 \times 10^{-5}.$$

Si los resultados de un experimento no tienen las mismas probabilidades de ocurrir, las probabilidades se deben asignar con base en el conocimiento previo o en la evidencia experimental. Por ejemplo, si una moneda no está balanceada, podemos estimar las probabilidades de caras y cruces lanzándola muchas veces y registrando los resultados. De acuerdo con la definición de **frecuencia relativa** de la probabilidad, las probabilidades verdaderas serían las fracciones de caras y cruces que ocurren a largo plazo. Otra forma intuitiva de comprender la probabilidad es el método de la **indiferencia**. Por ejemplo, si usted tiene un dado que cree que está balanceado, el método con el que podría determinar que hay $1/6$ de probabilidades de que resulte cada una de las seis caras después de lanzarlo una vez es el método de la indiferencia.

Para encontrar un valor numérico que represente de forma adecuada la probabilidad de ganar en el tenis, dependemos de nuestro desempeño previo en el juego, así como también del de nuestro oponente y, hasta cierto punto, de la capacidad de ganar que creemos tener. De manera similar, para calcular la probabilidad de que un caballo gane una carrera, debemos llegar a una probabilidad basada en las marcas anteriores de todos los caballos que participan en la carrera, así como de las marcas de los jinetes que los montan. La intuición, sin duda, también participa en la determinación del monto que estemos dispuestos a apostar. El uso de la intuición, las creencias personales y otra información indirecta para llegar a probabilidades se conoce como la definición **subjetiva** de la probabilidad.

En la mayoría de las aplicaciones de probabilidad de este libro la que opera es la interpretación de frecuencia relativa de probabilidad, la cual se basa en el experimento estadístico en vez de en la subjetividad y es considerada, más bien, como **frecuencia relativa limitante**. Como resultado, muchas aplicaciones de probabilidad en ciencia e ingeniería se deben basar en experimentos que se puedan repetir. Cuando asignamos probabilidades que se basan en información y opiniones previas, como en la afirmación: "hay grandes probabilidades de que los Gigantes pierdan el Súper Tazón", se encuentran

conceptos menos objetivos de probabilidad. Cuando las opiniones y la información previa difieren de un individuo a otro, la probabilidad subjetiva se vuelve el recurso pertinente. En la estadística bayesiana (véase el capítulo 18) se usará una interpretación más subjetiva de la probabilidad, la cual se basará en obtener información previa de probabilidad.

2.5 Reglas aditivas

A menudo resulta más sencillo calcular la probabilidad de algún evento a partir de las probabilidades conocidas de otros eventos. Esto puede ser cierto si el evento en cuestión se puede representar como la unión de otros dos eventos o como el complemento de algún evento. A continuación se presentan varias leyes importantes que con frecuencia simplifican el cálculo de las probabilidades. La primera, que se denomina **regla aditiva**, se aplica a uniones de eventos.

Teorema 2.7: Si A y B son dos eventos, entonces

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

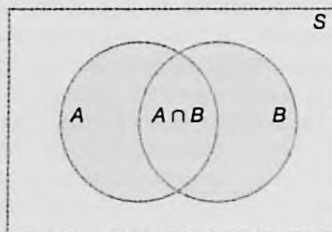


Figura 2.7: Regla aditiva de probabilidad.

Prueba: Considere el diagrama de Venn de la figura 2.7. $P(A \cup B)$ es la suma de las probabilidades de los puntos muestrales en $(A \cup B)$. Así, $P(A) + P(B)$ es la suma de todas las probabilidades en A más la suma de todas las probabilidades en B . Por lo tanto, sumamos dos veces las probabilidades en $(A \cap B)$. Como estas probabilidades se suman a $P(A \cap B)$, debemos restar esta probabilidad una vez para obtener la suma de las probabilidades en $A \cup B$. \square

Corolario 2.1: Si A y B son mutuamente excluyentes, entonces

$$P(A \cup B) = P(A) + P(B).$$

El corolario 2.1 es un resultado inmediato del teorema 2.7, pues si A y B son mutuamente excluyentes, $A \cap B = \emptyset$ y entonces $P(A \cap B) = P(\emptyset) = 0$. En general, podemos anotar el corolario 2.2.

Corolario 2.2: Si A_1, A_2, \dots, A_n son mutuamente excluyentes, entonces

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n).$$

Un conjunto de eventos $\{A_1, A_2, \dots, A_n\}$ de un espacio muestral S se denomina **partición** de S si A_1, A_2, \dots, A_n son mutuamente excluyentes y $A_1 \cup A_2 \cup \dots \cup A_n = S$. Por lo tanto, tenemos

Corolario 2.3: Si A_1, A_2, \dots, A_n es una partición de un espacio muestral S , entonces

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n) = P(S) = 1.$$

Como se esperaba, el teorema 2.7 se extiende de forma análoga.

Teorema 2.8: Para tres eventos A, B y C ,

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C).$$

Ejemplo 2.29: Al final del semestre John se va a graduar en la facultad de ingeniería industrial de una universidad. Después de tener entrevistas en dos empresas en donde quiere trabajar, determina que la probabilidad que tiene de lograr una oferta de empleo en la empresa A es 0.8, y que la probabilidad de obtenerla en la empresa B es 0.6. Si, por otro lado, considera que la probabilidad de recibir ofertas de ambas empresas es 0.5, ¿qué probabilidad tiene de obtener al menos una oferta de esas dos empresas?

Solución: Si usamos la regla aditiva tenemos

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = 0.8 + 0.6 - 0.5 = 0.9. \quad \blacksquare$$

Ejemplo 2.30: ¿Cuál es la probabilidad de obtener un total de 7 u 11 cuando se lanza un par de dados?

Solución: Sea A el evento de que resulte 7 y B el evento de que salga 11. Ahora bien, para 6 de los 36 puntos muestrales ocurre un total de 7 y sólo para 2 de ellos ocurre un total de 11. Como todos los puntos muestrales tienen la misma probabilidad, tenemos $P(A) = 1/6$ y $P(B) = 1/18$. Los eventos A y B son mutuamente excluyentes, ya que un total de 7 y uno de 11 no pueden ocurrir en el mismo lanzamiento. Por lo tanto,

$$P(A \cup B) = P(A) + P(B) = \frac{1}{6} + \frac{1}{18} = \frac{2}{9}.$$

Este resultado también se podría obtener contando el número total de puntos para el evento $A \cup B$, es decir, 8 y escribir

$$P(A \cup B) = \frac{n}{N} = \frac{8}{36} = \frac{2}{9}.$$

El teorema 2.7 y sus tres corolarios deberían ayudar al lector a comprender mejor la probabilidad y su interpretación. Los corolarios 2.1 y 2.2 sugieren el resultado muy intuitivo tratando con la probabilidad de que ocurra al menos uno de varios eventos, sin que puedan ocurrir dos de ellos simultáneamente. La probabilidad de que al menos ocurra uno es la suma de las probabilidades de ocurrencia de los eventos individuales. El tercer corolario simplemente establece que el valor mayor de una probabilidad (unidad) se asigna a todo el espacio muestral S .

Ejemplo 2.31: Las probabilidades de que un individuo que compra un automóvil nuevo elija uno de color verde, uno blanco, uno rojo o uno azul son 0.09, 0.15, 0.21 y 0.23, respectivamente. ¿cuál es la probabilidad de que un comprador dado adquiera un automóvil nuevo que tenga uno de esos colores?

Solución: Sean V , B , R y A los eventos de que un comprador seleccione, respectivamente, un automóvil verde, blanco, rojo o azul. Como estos cuatro eventos son mutuamente excluyentes, la probabilidad es

$$\begin{aligned} P(V \cup B \cup R \cup A) &= P(V) + P(B) + P(R) + P(A) \\ &= 0.09 + 0.15 + 0.21 + 0.23 = 0.68. \end{aligned} \quad \square$$

A menudo es más difícil calcular la probabilidad de que ocurra un evento que calcular la probabilidad de que el evento no ocurra. Si éste es el caso para algún evento A , simplemente calculamos primero $P(A')$ y, después, mediante el teorema 2.7, calculamos $P(A)$ por sustracción.

Teorema 2.9: Si A y A' son eventos complementarios, entonces

$$P(A) + P(A') = 1$$

Prueba: Como $A \cup A' = S$, y los conjuntos A y A' son disjuntos, entonces

$$1 = P(S) = P(A \cup A') = P(A) + P(A') \quad \square$$

Ejemplo 2.32: Si las probabilidades de que un mecánico automotriz dé servicio a 3, 4, 5, 6, 7, 8 o más vehículos en un día de trabajo dado son 0.12, 0.19, 0.28, 0.24, 0.10 y 0.07, respectivamente, ¿cuál es la probabilidad de que dé servicio al menos a 5 vehículos el siguiente día de trabajo?

Solución: Sea E el evento de que al menos 5 automóviles reciban servicio. Ahora bien, $P(E) = 1 - P(E')$, donde E' es el evento de que menos de 5 automóviles reciban servicio. Como

$$P(E') = 0.12 + 0.19 = 0.31.$$

del teorema 2.9 se deduce que

$$P(E) = 1 - 0.31 = 0.69. \quad \square$$

Ejemplo 2.33: Suponga que las especificaciones del fabricante para la longitud del cable de cierto tipo de computadora son 2000 ± 10 milímetros. En esta industria se sabe que el cable pequeño tiene la misma probabilidad de salir defectuoso (de no cumplir con las especificaciones) que el cable grande. Es decir, la probabilidad de que aleatoriamente se produzca un

cable con una longitud mayor que 2010 milímetros es igual a la probabilidad de producirlo con una longitud menor que 1990 milímetros. Se sabe que la probabilidad de que el procedimiento de producción cumpla con las especificaciones es 0.99.

- a) ¿Cuál es la probabilidad de que un cable elegido al azar sea muy largo?
 b) ¿Cuál es la probabilidad de que un cable elegido al azar sea más grande que 1990 milímetros?

Solución: Sea E el evento de que un cable cumpla con las especificaciones. Sean P y G los eventos de que el cable sea muy pequeño o muy grande, respectivamente. Entonces,

- a) $P(E) = 0.99$ y $P(P) = P(G) = (1 - 0.99)/2 = 0.005$.
 b) Si la longitud de un cable seleccionado al azar se denota con X , tenemos

$$P(1990 \leq X \leq 2010) = P(E) = 0.99.$$

$$\text{Como } P(X \geq 2010) = P(G) = 0.005,$$

$$P(X \geq 1990) = P(E) + P(G) = 0.995$$

Esto también se resuelve utilizando el teorema 2.9:

$$P(X \geq 1990) + P(X < 1990) = 1.$$

$$\text{Así, } P(X \geq 1990) = 1 - P(P) = 1 - 0.005 = 0.995.$$

Ejercicios

2.49 Encuentre los errores en cada una de las siguientes aseveraciones:

- a) Las probabilidades de que un vendedor de automóviles venda 0, 1, 2 o 3 unidades en un día dado de febrero son 0.19, 0.38, 0.29 y 0.15, respectivamente.
 b) La probabilidad de que llueva mañana es 0.40 y la probabilidad de que no llueva es 0.52.
 c) Las probabilidades de que una impresora cometa 0, 1, 2, 3 o 4 o más errores al imprimir un documento son 0.19, 0.34, -0.25, 0.43 y 0.29, respectivamente.
 d) Al sacar una carta de una baraja en un solo intento la probabilidad de seleccionar un corazón es $1/4$, la probabilidad de seleccionar una carta negra es $1/2$, y la probabilidad de seleccionar una carta de corazones y negra es $1/8$.

2.50 Suponga que todos los elementos de S en el ejercicio 2.8 de la página 42 tienen la misma probabilidad de ocurrencia y calcule

- a) la probabilidad del evento A ;
 b) la probabilidad del evento C ;
 c) la probabilidad del evento $A \cap C$.

2.51 Una caja contiene 500 sobres, de los cuales 75 contienen \$100 en efectivo, 150 contienen \$25 y 275 contienen \$10. Se puede comprar un sobre en \$25. ¿Cuál es el espacio muestral para las diferentes cantidades de dinero? Asigne probabilidades a los puntos muestrales y después calcule la probabilidad de que el primer sobre que se compre contenga menos de \$100.

2.52 Suponga que se descubre que, en un grupo de 500 estudiantes universitarios de último año, 210 fuman, 258 consumen bebidas alcohólicas, 216 comen entre comidas, 122 fuman y consumen bebidas alcohólicas, 83 comen entre comidas y consumen bebidas alcohólicas, 97 fuman y comen entre comidas y 52 tienen esos tres hábitos nocivos para la salud. Si se selecciona al azar a un miembro de este grupo, calcule la probabilidad de que el estudiante

- a) fume pero no consuma bebidas alcohólicas;
 b) coma entre comidas y consuma bebidas alcohólicas pero no fume;
 c) no fume ni coma entre comidas.

2.53 La probabilidad de que una industria estadounidense se ubique en Shanghái, China, es 0.7, la probabilidad de que se ubique en Beijing, China, es 0.4 y la

probabilidad de que se ubique en Shanghai o Beijing, o en ambas ciudades, es 0.8. ¿Cuál es la probabilidad de que la industria se ubique...

- en ambas ciudades?
- en ninguna de esas ciudades?

2.54 Basado en su experiencia, un agente bursátil considera que en las condiciones económicas actuales la probabilidad de que un cliente invierta en bonos libres de impuestos es 0.6, la de que invierta en fondos comunes de inversión es 0.3 y la de que invierta en ambos es 0.15. En esta ocasión encuentre la probabilidad de que un cliente invierta

- en bonos libres de impuestos o en fondos comunes de inversión;
- en ninguno de esos dos instrumentos.

2.55 Si cada artículo codificado en un catálogo empieza con 3 letras distintas seguidas por 4 dígitos distintos de cero, calcule la probabilidad de seleccionar aleatoriamente uno de estos artículos codificados que tenga como primera letra una vocal y el último dígito sea par.

2.56 Un fabricante de automóviles está preocupado por el posible retiro de su sedán de cuatro puertas con mayor venta. Si fuera retirado habría 0.25 de probabilidad de que haya un defecto en el sistema de frenos, 0.18 de que haya un defecto en la transmisión, 0.17 de que esté en el sistema de combustible y 0.40 de que esté en alguna otra área.

- ¿Cuál es la probabilidad de que el defecto esté en los frenos o en el sistema de combustible, si la probabilidad de que haya defectos en ambos sistemas de manera simultánea es 0.15?
- ¿Cuál es la probabilidad de que no haya defecto en los frenos o en el sistema de combustible?

2.57 Si se elige al azar una letra del alfabeto inglés, encuentre la probabilidad de que la letra

- sea una vocal excepto y;
- esté listada en algún lugar antes de la letra j;
- esté listada en algún lugar después de la letra g.

2.58 Se lanza un par de dados. Calcule la probabilidad de obtener

- un total de 8;
- máximo un total de 5.

2.59 En una mano de póquer que consta de 5 cartas, encuentre la probabilidad de tener

- 3 ases;
- 4 cartas de corazones y 1 de tréboles.

2.60 Si se toman 3 libros al azar, de un librero que contiene 5 novelas, 3 libros de poemas y 1 diccionario, ¿cuál es la probabilidad de que...

- se seleccione el diccionario?
- se seleccionen 2 novelas y 1 libro de poemas?

2.61 En un grupo de 100 estudiantes graduados de preparatoria, 54 estudiaron matemáticas, 69 estudiaron historia y 35 cursaron matemáticas e historia. Si se selecciona al azar uno de estos estudiantes, calcule la probabilidad de que

- el estudiante haya cursado matemáticas o historia;
- el estudiante no haya llevado ninguna de estas materias;
- el estudiante haya cursado historia pero no matemáticas.

2.62 La empresa Dom's Pizza utiliza pruebas de sabor y el análisis estadístico de los datos antes de comercializar cualquier producto nuevo. Considere un estudio que incluye tres tipos de pastas (delgada, delgada con ajo y orégano, y delgada con trozos de queso). Dom's también está estudiando tres salsas (estándar, una nueva salsa con más ajo y una nueva salsa con albahaca fresca).

- ¿Cuántas combinaciones de pasta y salsa se incluyen?
- ¿Cuál es la probabilidad de que un juez reciba una pasta delgada sencilla con salsa estándar en su primera prueba de sabor?

2.63 A continuación se listan los porcentajes, proporcionados por *Consumer Digest* (julio/agosto de 1996), de las probables ubicaciones de las PC en una casa:

Dormitorio de adultos:	0.03
Dormitorio de niños:	0.15
Otro dormitorio:	0.14
Oficina o estudio:	0.40
Otra habitación:	0.28

- ¿Cuál es la probabilidad de que una PC esté en un dormitorio?
- ¿Cuál es la probabilidad de que no esté en un dormitorio?
- Suponga que de entre las casas que tienen una PC se selecciona una al azar, ¿en qué habitación esperaría encontrar una PC?

2.64 Existe interés por la vida de un componente electrónico. Suponga que se sabe que la probabilidad de que el componente funcione más de 6000 horas es 0.42. Suponga, además, que la probabilidad de que el componente *no dure más de 4000 horas* es 0.04.

- ¿Cuál es la probabilidad de que la vida del componente sea menor o igual a 6000 horas?
- ¿Cuál es la probabilidad de que la vida del componente sea mayor que 4000 horas?

2.65 Considere la situación del ejercicio 2.64. Sea A el evento de que el componente falle en una prueba específica y B el evento de que se deforme pero no falle. El evento A ocurre con una probabilidad de 0.20 y el evento B ocurre con una probabilidad de 0.35.

- a) ¿Cuál es la probabilidad de que el componente no falle en la prueba?
- b) ¿Cuál es la probabilidad de que el componente funcione perfectamente bien (es decir, que ni se deforme ni falle en la prueba)?
- c) ¿Cuál es la probabilidad de que el componente falle o se deforme en la prueba?

2.66 A los obreros de las fábricas se les motiva constantemente a practicar la tolerancia cero para prevenir accidentes en el lugar de trabajo. Los accidentes pueden ocurrir porque el ambiente o las condiciones laborales son inseguros. Por otro lado, los accidentes pueden ocurrir por negligencia o fallas humanas. Además, los horarios de trabajo de 7:00 A.M. a 3:00 P.M. (turno matutino), de 3:00 P.M. a 11:00 P.M. (turno vespertino) y de 11:00 P.M. a 7:00 A.M. (turno nocturno) podría ser un factor. El año pasado ocurrieron 300 accidentes. Los porcentajes de los accidentes por la combinación de condiciones son los que siguen:

Turno	Condiciones inseguras	Fallas humanas
Matutino	5%	32%
Vespertino	6%	25%
Nocturno	2%	30%

Si se elige aleatoriamente un reporte de accidente de entre los 300 reportes,

- a) ¿Cuál es la probabilidad de que el accidente haya ocurrido en el turno nocturno?
- b) ¿Cuál es la probabilidad de que el accidente haya ocurrido debido a una falla humana?
- c) ¿Cuál es la probabilidad de que el accidente haya ocurrido debido a las condiciones inseguras?
- d) ¿Cuál es la probabilidad de que el accidente haya ocurrido durante los turnos vespertino o nocturno?

2.67 Considere la situación del ejemplo 2.32 de la página 58.

- a) ¿Cuál es la probabilidad de que el número de automóviles que recibirán servicio del mecánico no sea mayor de 4?
- b) ¿Cuál es la probabilidad de que el mecánico dé servicio a menos de 8 automóviles?
- c) ¿Cuál es la probabilidad de que el mecánico dé servicio a 3 o 4 automóviles?

2.68 Existe interés por el tipo de horno, eléctrico o de gas, que se compra en una tienda departamental específica. Considere la decisión que al respecto toman seis clientes distintos.

- a) Suponga que hay 0.40 de probabilidades de que como máximo dos de esos clientes compren un horno eléctrico. ¿Cuál será la probabilidad de que al menos tres compren un horno eléctrico?

- b) Suponga que se sabe que la probabilidad de que los seis compren el horno eléctrico es 0.007, mientras que la probabilidad de que los seis compren el horno de gas es 0.104. ¿Cuál es la probabilidad de vender, por lo menos, un horno de cada tipo?

2.69 En muchas áreas industriales es común que se utilicen máquinas para llenar las cajas de productos. Esto ocurre tanto en la industria de comestibles como en otras que fabrican productos de uso doméstico, como los detergentes. Dichas máquinas no son perfectas y, de hecho, podrían cumplir las especificaciones de llenado de las cajas (A), llenarlas por debajo del nivel especificado (B) o rebasar el límite de llenado (C). Por lo general, lo que se busca evitar es la práctica del llenado insuficiente. Sea $P(B) = 0.001$, mientras que $P(A) = 0.990$.

- a) Determine $P(C)$.
- b) ¿Cuál es la probabilidad de que la máquina no llene de manera suficiente?
- c) ¿Cuál es la probabilidad de que la máquina llene de más o de menos?

2.70 Considere la situación del ejercicio 2.69. Suponga que se producen 50,000 cajas de detergente por semana, y que los clientes "devuelven" las cajas que no están suficientemente llenas y solicitan que se les reembolse lo que pagaron por ellas. Suponga que se sabe que el "costo" de producción de cada caja es de \$4.00 y que se venden a \$4.50.

- a) ¿Cuál es la utilidad semanal cuando no hay devoluciones de cajas defectuosas?
- b) ¿Cuál es la pérdida en utilidades esperada debido a la devolución de cajas insuficientemente llenadas?

2.71 Como podría sugerir la situación del ejercicio 2.69, a menudo los procedimientos estadísticos se utilizan para control de calidad (es decir, control de calidad industrial). A veces el *peso* de un producto es una variable importante que hay que controlar. Se dan especificaciones de peso para ciertos productos empacados, y si un paquete no las cumple (está muy ligero o muy pesado) se rechaza. Los datos históricos sugieren que la probabilidad de que un producto empacado cumpla con las especificaciones de peso es 0.95; mientras que la probabilidad de que sea demasiado ligero es 0.002. El fabricante invierte \$20.00 en la producción de cada uno de los productos empacados y el consumidor los adquiere a un precio de \$25.00.

- a) ¿Cuál es la probabilidad de que un paquete elegido al azar de la línea de producción sea demasiado pesado?
- b) Si todos los paquetes cumplen con las especificaciones de peso, ¿qué utilidad recibirá el fabricante por cada 10,000 paquetes que venda?

- c) Suponga que todos los paquetes defectuosos fueron rechazados y perdieron todo su valor, ¿a cuánto se reduciría la utilidad de la venta de 10,000 paquetes debido a que no se cumplieron las especificaciones de peso?

2.72 Demuestre que

$$P(A' \cap B') = 1 + P(A \cap B) - P(A) - P(B).$$

2.6 Probabilidad condicional, independencia y regla del producto

Un concepto muy importante en la teoría de probabilidad es la probabilidad condicional. En algunas aplicaciones el profesional se interesa por la estructura de probabilidad bajo ciertas restricciones. Por ejemplo, en epidemiología, en lugar de estudiar las probabilidades de que una persona de la población general tenga diabetes, podría ser más interesante conocer esta probabilidad en un grupo distinto, como el de las mujeres asiáticas cuya edad está en el rango de 35 a 50 años, o como el de los hombres hispanos cuya edad está entre los 40 y los 60 años. A este tipo de probabilidad se le conoce como probabilidad condicional.

Probabilidad condicional

La probabilidad de que ocurra un evento B cuando se sabe que ya ocurrió algún evento A se llama **probabilidad condicional** y se denota con $P(B|A)$. El símbolo $P(B|A)$ por lo general se lee como “la probabilidad de que ocurra B , dado que ocurrió A ”, o simplemente, “la probabilidad de B , dado A ”.

Considere el evento B de obtener un cuadrado perfecto cuando se lanza un dado. El dado se construye de modo que los números pares tengan el doble de probabilidad de ocurrencia que los números nones. Con base en el espacio muestral $S = \{1, 2, 3, 4, 5, 6\}$, en el que a los números impares y a los pares se les asignaron probabilidades de $1/9$ y $2/9$, respectivamente, la probabilidad de que ocurra B es de $1/3$. Suponga ahora que se sabe que el lanzamiento del dado tiene como resultado un número mayor que 3. Tenemos ahora un espacio muestral reducido, $A = \{4, 5, 6\}$, que es un subconjunto de S . Para encontrar la probabilidad de que ocurra B , en relación con el espacio muestral A , debemos comenzar por asignar nuevas probabilidades a los elementos de A , que sean proporcionales a sus probabilidades originales de modo que su suma sea 1. Al asignar una probabilidad de w al número non en A y una probabilidad de $2w$ a los dos números pares, tenemos $5w = 1$ o $w = 1/5$. En relación con el espacio A , encontramos que B contiene sólo el elemento 4. Si denotamos este evento con el símbolo $B|A$, escribimos $B|A = \{4\}$ y, en consecuencia,

$$P(B|A) = \frac{2}{5}.$$

Este ejemplo ilustra que los eventos pueden tener probabilidades diferentes cuando se consideran en relación con diferentes espacios muestrales.

También podemos escribir

$$P(B|A) = \frac{2}{5} = \frac{2/9}{5/9} = \frac{P(A \cap B)}{P(A)},$$

donde $P(A \cap B)$ y $P(A)$ se calculan a partir del espacio muestral original S . En otras palabras, una probabilidad condicional relativa a un subespacio A de S se puede calcular en forma directa de las probabilidades que se asignan a los elementos del espacio muestral original S .

Definición 2.10: La probabilidad condicional de B , dado A , que se denota con $P(B|A)$, se define como

$$P(B|A) = \frac{P(A \cap B)}{P(A)}, \text{ siempre que } P(A) > 0.$$

Un ejemplo más: suponga que tenemos un espacio muestral S constituido por la población de adultos de una pequeña ciudad que cumplen con los requisitos para obtener un título universitario. Debemos clasificarlos de acuerdo con su género y situación laboral. Los datos se presentan en la tabla 2.1.

Tabla 2.1: Clasificación de los adultos de una pequeña ciudad

	Empleado	Desempleado	Total
Hombre	460	40	500
Mujer	140	260	400
Total	600	300	900

Se seleccionará al azar a uno de estos individuos para que realice un viaje a través del país con el fin de promover las ventajas de establecer industrias nuevas en la ciudad. Nos interesaremos en los eventos siguientes:

M : se elige a un hombre,

E : el elegido tiene empleo.

Al utilizar el espacio muestral reducido E , encontramos que

$$P(M|E) = \frac{460}{600} = \frac{23}{30}.$$

Sea $n(A)$ el número de elementos en cualquier conjunto A . Podemos utilizar esta notación, puesto que cada uno de los adultos tiene las mismas probabilidades de ser elegido, para escribir

$$P(M|E) = \frac{n(E \cap M)}{n(E)} = \frac{n(E \cap M)/n(S)}{n(E)/n(S)} = \frac{P(E \cap M)}{P(E)},$$

en donde $P(E \cap M)$ y $P(E)$ se calculan a partir del espacio muestral original S . Para verificar este resultado observe que

$$P(E) = \frac{600}{900} = \frac{2}{3} \text{ y } P(E \cap M) = \frac{460}{900} = \frac{23}{45}.$$

Por lo tanto,

$$P(M|E) = \frac{23/45}{2/3} = \frac{23}{30},$$

como antes.

Ejemplo 2.34: La probabilidad de que un vuelo programado normalmente salga a tiempo es $P(D) = 0.83$, la probabilidad de que llegue a tiempo es $P(A) = 0.82$ y la probabilidad de que

salga y llegue a tiempo es $P(D \cap A) = 0.78$. Calcule la probabilidad de que un avión a) llegue a tiempo, dado que salió a tiempo; y b) salió a tiempo, dado que llegó a tiempo.

Solución: Al utilizar la definición 2.10 tenemos lo que sigue:

a) La probabilidad de que un avión llegue a tiempo, dado que salió a tiempo es

$$P(A|D) = \frac{P(D \cap A)}{P(D)} = \frac{0.78}{0.83} = 0.94.$$

b) La probabilidad de que un avión haya salido a tiempo, dado que llegó a tiempo es

$$P(D|A) = \frac{P(D \cap A)}{P(A)} = \frac{0.78}{0.82} = 0.95. \quad \blacksquare$$

La noción de probabilidad condicional brinda la capacidad de reevaluar la idea de probabilidad de un evento a la luz de la información adicional; es decir, cuando se sabe que ocurrió otro evento. La probabilidad $P(A|B)$ es una actualización de $P(A)$ basada en el conocimiento de que ocurrió el evento B . En el ejemplo 2.34 es importante conocer la probabilidad de que el vuelo llegue a tiempo. Tenemos la información de que el vuelo no salió a tiempo. Con esta información adicional, la probabilidad más pertinente es $P(A|D')$, esto es, la probabilidad de que llegue a tiempo, dado que no salió a tiempo. A menudo las conclusiones que se obtienen a partir de observar la probabilidad condicional más importante cambian drásticamente la situación. En este ejemplo, el cálculo de $P(A|D')$ es

$$P(A|D') = \frac{P(A \cap D')}{P(D')} = \frac{0.82 - 0.78}{0.17} = 0.24.$$

Como resultado, la probabilidad de una llegada a tiempo disminuye significativamente ante la presencia de la información adicional.

Ejemplo 2.35: El concepto de probabilidad condicional tiene innumerables aplicaciones industriales y biomédicas. Considere un proceso industrial en el ramo textil, en el que se producen listones de una tela específica. Los listones pueden resultar con defectos en dos de sus características: la longitud y la textura. En el segundo caso el proceso de identificación es muy complicado. A partir de información histórica del proceso se sabe que 10% de los listones no pasan la prueba de longitud, que 5% no pasan la prueba de textura y que sólo 0.8% no pasan ninguna de las dos pruebas. Si en el proceso se elige un listón al azar y una medición rápida identifica que no pasa la prueba de longitud, ¿cuál es la probabilidad de que la textura esté defectuosa?

Solución: Considere los eventos

L : defecto en longitud,

T : defecto en textura.

Dado que el listón tiene una longitud defectuosa, la probabilidad de que este listón tenga una textura defectuosa está dada por

$$P(T|L) = \frac{P(T \cap L)}{P(L)} = \frac{0.008}{0.1} = 0.08.$$

Eventos independientes

En el experimento del lanzamiento de un dado de la página 62 señalamos que $P(B|A) = 2/5$, mientras que $P(B) = 1/3$. Es decir, $P(B|A) \neq P(B)$, lo cual indica que B depende de A . Consideremos ahora un experimento en el que se sacan 2 cartas, una después de la otra, de una baraja ordinaria, con reemplazo. Los eventos se definen como

A : la primera carta es un as,

B : la segunda carta es una espada.

Como la primera carta se reemplaza, nuestro espacio muestral para la primera y segunda cartas consta de 52 cartas, que contienen 4 ases y 13 espadas. Entonces,

$$P(B|A) = \frac{13}{52} = \frac{1}{4} \quad \text{y} \quad P(B) = \frac{13}{52} = \frac{1}{4}.$$

Es decir, $P(B|A) = P(B)$. Cuando esto es cierto, se dice que los eventos A y B son **independientes**.

Aunque la probabilidad condicional permite alterar la probabilidad de un evento a la luz de material adicional, también nos permite entender mejor el muy importante concepto de **independencia** o, en el contexto actual, de eventos independientes. En el ejemplo 2.34 del aeropuerto, $P(A|D)$ difiere de $P(A)$. Esto sugiere que la ocurrencia de D influye en A y esto es lo que, de hecho, se espera en este caso. Sin embargo, considere la situación en donde tenemos los eventos A y B , y

$$P(A|B) = P(A).$$

En otras palabras, la ocurrencia de B no influye en las probabilidades de ocurrencia de A . Aquí la ocurrencia de A es independiente de la ocurrencia de B . No podemos dejar de resaltar la importancia del concepto de independencia, ya que desempeña un papel vital en el material de casi todos los capítulos de este libro y en todas las áreas de la estadística aplicada.

Definición 2.11: Dos eventos A y B son **independientes** si y sólo si

$$P(B|A) = P(B) \quad \text{o} \quad P(A|B) = P(A),$$

si se asume la existencia de probabilidad condicional. De otra forma, A y B son **dependientes**.

La condición $P(B|A) = P(B)$ implica que $P(A|B) = P(A)$, y viceversa. Para los experimentos de extracción de una carta, donde mostramos que $P(B|A) = P(B) = 1/4$, también podemos ver que $P(A|B) = P(A) = 1/13$.

La regla de producto o regla multiplicativa

Al multiplicar la fórmula de la definición 2.10 por $P(A)$, obtenemos la siguiente **regla multiplicativa** importante (o **regla de producto**), que nos permite calcular la probabilidad de que ocurran dos eventos.

Teorema 2.10: Si en un experimento pueden ocurrir los eventos A y B , entonces

$$P(A \cap B) = P(A)P(B|A), \text{ siempre que } P(A) > 0.$$

Por consiguiente, la probabilidad de que ocurran A y B es igual a la probabilidad de que ocurra A multiplicada por la probabilidad condicional de que ocurra B , dado que ocurre A . Como los eventos $A \cap B$ y $B \cap A$ son equivalentes, del teorema 2.10 se deduce que también podemos escribir

$$P(A \cap B) = P(B \cap A) = P(B)P(A|B).$$

En otras palabras, no importa qué evento se considere como A ni qué evento se considere como B .

Ejemplo 2.36: Suponga que tenemos una caja de fusibles que contiene 20 unidades, de las cuales 5 están defectuosas. Si se seleccionan 2 fusibles al azar y se retiran de la caja, uno después del otro, sin reemplazar el primero, ¿cuál es la probabilidad de que ambos fusibles estén defectuosos?

Solución: Sean A el evento de que el primer fusible esté defectuoso y B el evento de que el segundo esté defectuoso; entonces, interpretamos $A \cap B$ como el evento de que ocurra A , y entonces B ocurre después de que haya ocurrido A . La probabilidad de sacar primero un fusible defectuoso es $1/4$; entonces, la probabilidad de separar un segundo fusible defectuoso de los restantes 4 es $4/19$. Por lo tanto,

$$P(A \cap B) = \left(\frac{1}{4}\right) \left(\frac{4}{19}\right) = \frac{1}{19}.$$

Ejemplo 2.37: Una bolsa contiene 4 bolas blancas y 3 negras, y una segunda bolsa contiene 3 blancas y 5 negras. Se saca una bola de la primera bolsa y se coloca sin verla en la segunda bolsa. ¿Cuál es la probabilidad de que ahora se saque una bola negra de la segunda bolsa?

Solución: N_1 , N_2 y B_1 representan, respectivamente, la extracción de una bola negra de la bolsa 1, una bola negra de la bolsa 2 y una bola blanca de la bolsa 1. Nos interesa la unión de los eventos mutuamente excluyentes $N_1 \cap N_2$ y $B_1 \cap N_2$. Las diversas posibilidades y sus probabilidades se ilustran en la figura 2.8. Entonces

$$\begin{aligned} P[(N_1 \cap N_2) \cup (B_1 \cap N_2)] &= P(N_1 \cap N_2) + P(B_1 \cap N_2) \\ &= P(N_1)P(N_2|N_1) + P(B_1)P(N_2|B_1) \\ &= \left(\frac{3}{7}\right) \left(\frac{6}{9}\right) + \left(\frac{4}{7}\right) \left(\frac{5}{9}\right) = \frac{38}{63}. \end{aligned}$$

Si, en el ejemplo 2.36, el primer fusible se reemplaza y los fusibles se acomodan por completo antes de extraer el segundo, entonces la probabilidad de que se extraiga un fusible defectuoso en la segunda selección sigue siendo $1/4$; es decir, $P(B|A) = P(B)$, y los eventos A y B son independientes. Cuando esto es cierto podemos sustituir $P(B)$ por $P(B|A)$ en el teorema 2.10 para obtener la siguiente regla multiplicativa especial.

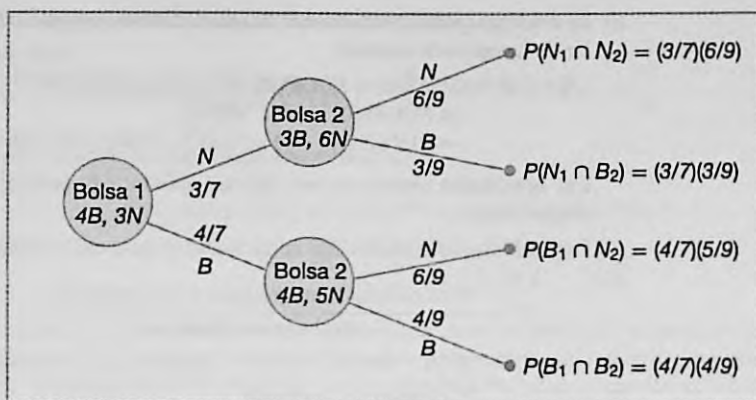


Figura 2.8: Diagrama de árbol para el ejemplo 2.37.

Teorema 2.11: Dos eventos A y B son independientes si y sólo si

$$P(A \cap B) = P(A)P(B).$$

Por lo tanto, para obtener la probabilidad de que ocurran dos eventos independientes simplemente calculamos el producto de sus probabilidades individuales.

Ejemplo 2.38: Una pequeña ciudad dispone de un carro de bomberos y una ambulancia para emergencias. La probabilidad de que el carro de bomberos esté disponible cuando se necesite es 0.98 y la probabilidad de que la ambulancia esté disponible cuando se la requiera es 0.92. En el evento de un herido en un incendio, calcule la probabilidad de que tanto la ambulancia como el carro de bomberos estén disponibles, suponiendo que operan de forma independiente.

Solución: Sean A y B los respectivos eventos de que estén disponibles el carro de bomberos y la ambulancia. Entonces,

$$P(A \cap B) = P(A)P(B) = (0.98)(0.92) = 0.9016. \quad \blacksquare$$

Ejemplo 2.39: Un sistema eléctrico consta de cuatro componentes, como se ilustra en la figura 2.9. El sistema funciona si los componentes A y B funcionan, y si funciona cualquiera de los componentes C o D . La confiabilidad (probabilidad de que funcionen) de cada uno de los componentes también se muestra en la figura 2.9. Calcule la probabilidad de a) que el sistema completo funcione y de b) que el componente C no funcione, dado que el sistema completo funciona. Suponga que los cuatro componentes funcionan de manera independiente.

Solución: En esta configuración del sistema, A , B y el subsistema C y D constituyen un sistema de circuitos en serie; mientras que el subsistema C y D es un sistema de circuitos en paralelo.

- a) Es evidente que la probabilidad de que el sistema completo funcione se puede calcular de la siguiente manera:

$$\begin{aligned} P[A \cap B \cap (C \cup D)] &= P(A)P(B)P(C \cup D) = P(A)P(B)[1 - P(C' \cap D')] \\ &= P(A)P(B)[1 - P(C')P(D')] \\ &= (0.9)(0.9)[1 - (1 - 0.8)(1 - 0.8)] = 0.7776. \end{aligned}$$

Las igualdades anteriores son válidas debido a la independencia entre los cuatro componentes.

- b) Para calcular la probabilidad condicional en este caso, observe que

$$\begin{aligned} P &= \frac{P(\text{el sistema funciona pero } C \text{ no funciona})}{P(\text{el sistema funciona})} \\ &= \frac{P(A \cap B \cap C' \cap D)}{P(\text{el sistema funciona})} = \frac{(0.9)(0.9)(1 - 0.8)(0.8)}{0.7776} = 0.1667. \end{aligned}$$

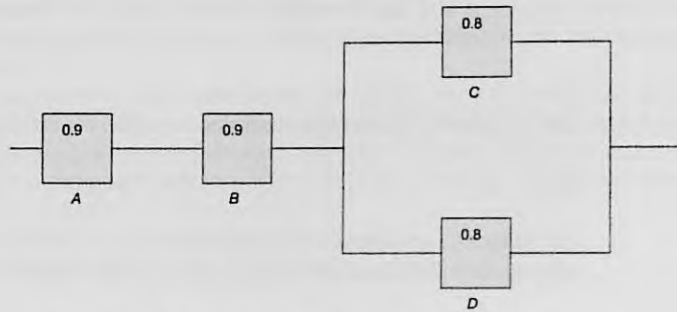


Figura 2.9: Un sistema eléctrico para el ejemplo 2.39.

La regla multiplicativa se puede extender a situaciones con más de dos eventos.

Teorema 2.12: Si, en un experimento, pueden ocurrir los eventos A_1, A_2, \dots, A_k , entonces

$$\begin{aligned} P(A_1 \cap A_2 \cap \dots \cap A_k) \\ &= P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2) \dots P(A_k|A_1 \cap A_2 \cap \dots \cap A_{k-1}). \end{aligned}$$

Si los eventos A_1, A_2, \dots, A_k son independientes, entonces

$$P(A_1 \cap A_2 \cap \dots \cap A_k) = P(A_1)P(A_2) \dots P(A_k)$$

Ejemplo 2.40: Se sacan tres cartas seguidas, sin reemplazo, de una baraja ordinaria. Encuentre la probabilidad de que ocurra el evento $A_1 \cap A_2 \cap A_3$, donde A_1 es el evento de que la primera carta sea un as rojo, A_2 el evento de que la segunda carta sea un 10 o una jota y A_3 el evento de que la tercera carta sea mayor que 3 pero menor que 7.

Solución: Primero definimos los eventos:

A_1 : la primera carta es un as rojo,

A_2 : la segunda carta es un 10 o una jota,

A_3 : la tercera carta es mayor que 3 pero menor que 7.

Ahora bien,

$$P(A_1) = \frac{2}{52}, P(A_2|A_1) = \frac{8}{51}, P(A_3|A_1 \cap A_2) = \frac{12}{50},$$

por lo tanto, por medio del teorema 2.12,

$$\begin{aligned} P(A_1 \cap A_2 \cap A_3) &= P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2) \\ &= \left(\frac{2}{52}\right) \left(\frac{8}{51}\right) \left(\frac{12}{50}\right) = \frac{8}{5525}. \end{aligned}$$

La propiedad de independencia establecida en el teorema 2.11 se puede extender a situaciones con más de dos eventos. Considere, por ejemplo, el caso de los tres eventos A , B y C . No basta con tener $P(A \cap B \cap C) = P(A)P(B)P(C)$ como una definición de independencia entre los tres. Suponga que $A = B \cap C = \emptyset$, el conjunto vacío. Aunque $A \cap B \cap C = \emptyset$, que da como resultado $P(A \cap B \cap C) = 0 = P(A)P(B)P(C)$, los eventos A y B no son independientes. En consecuencia, tenemos la siguiente definición:

Definición 2.12: Un conjunto de eventos $\mathcal{A} = \{A_1, \dots, A_n\}$ son mutuamente independientes si para cualquier subconjunto de \mathcal{A} , A_{i_1}, \dots, A_{i_k} , para $k \leq n$, tenemos

$$P(A_{i_1} \cap \dots \cap A_{i_k}) = P(A_{i_1}) \dots P(A_{i_k}).$$

Ejercicios

2.73 Si R es el evento de que un convicto cometa un robo a mano armada y D es el evento de que el convicto venda drogas, exprese en palabras lo que en probabilidades se indica como

- $P(R|D)$;
- $P(D|R)$;
- $P(R|D')$.

2.74 Un grupo de estudiantes de física avanzada se compone de 10 alumnos de primer año, 30 del último año y 10 graduados. Las calificaciones finales muestran que 3 estudiantes de primer año, 10 del último año y 5 de los graduados obtuvieron 10 en el curso. Si se elige un estudiante al azar de este grupo y se descubre que es uno de los que obtuvieron 10 de calificación, ¿cuál es la probabilidad de que sea un estudiante de último año?

2.75 La siguiente es una clasificación, según el género y el nivel de escolaridad, de una muestra aleatoria de 200 adultos.

Escolaridad	Hombre	Mujer
Primaria	38	45
Secundaria	28	50
Universidad	22	17

Si se elige una persona al azar de este grupo, ¿cuál es la probabilidad de que...

- la persona sea hombre, dado que su escolaridad es de secundaria?;
- la persona no tenga un grado universitario, dado que es mujer?

2.76 En un experimento para estudiar la relación que existe entre el hábito de fumar y la hipertensión arterial se reúnen los siguientes datos para 180 individuos:

	No fumadores	Fumadores moderados	Fumadores empedernidos
H	21	36	30
SH	48	26	19

donde las letras H y SH de la tabla representan *Hipertensión* y *Sin hipertensión*, respectivamente. Si se selecciona uno de estos individuos al azar, calcule la probabilidad de que la persona...

- sufra hipertensión, dado que es una fumadora empedernida;
- no fume, dado que no padece hipertensión.

2.77 En un grupo de 100 estudiantes de bachillerato que están cursando el último año, 42 cursaron matemáticas, 68 psicología, 54 historia, 22 matemáticas e historia, 25 matemáticas y psicología, 7 historia pero ni matemáticas ni psicología, 10 las tres materias y 8 no cursaron ninguna de las tres. Seleccione al azar a un

estudiante de este grupo y calcule la probabilidad de los siguientes eventos:

- Una persona inscrita en psicología y cursa las tres materias;
- Una persona que no está inscrita en psicología y esté cursando historia y matemáticas.

2.78 Un fabricante de una vacuna para la gripe está interesado en determinar la calidad de su suero. Con ese fin tres departamentos diferentes procesan los lotes de suero y tienen tasas de rechazo de 0.10, 0.08 y 0.12, respectivamente. Las inspecciones de los tres departamentos son secuenciales e independientes.

- ¿Cuál es la probabilidad de que un lote de suero sobreviva a la primera inspección departamental pero sea rechazado por el segundo departamento?
- ¿Cuál es la probabilidad de que un lote de suero sea rechazado por el tercer departamento?

2.79 En *USA Today* (5 de septiembre de 1996) se listaron los siguientes resultados de una encuesta sobre el uso de ropa para dormir mientras se viaja:

	Hombre	Mujer	Total
Ropa interior	0.020	0.024	0.244
Camisón	0.002	0.180	0.182
Nada	0.160	0.018	0.178
Pijama	0.102	0.073	0.175
Camiseta	0.046	0.088	0.134
Otros	0.084	0.003	0.087

- ¿Cuál es la probabilidad de que un viajero sea una mujer que duerme desnuda?
- ¿Cuál es la probabilidad de que un viajero sea hombre?
- Si el viajero fuera hombre, ¿cuál sería la probabilidad de que duerma con pijama?
- ¿Cuál es la probabilidad de que un viajero sea hombre si duerme con pijama o con camiseta?

2.80 La probabilidad de que cuando se tenga que llenar el tanque de gasolina de un automóvil también se necesite cambiarle el aceite es 0.25, la probabilidad de que también se le tenga que cambiar el filtro de aceite es 0.40, y la probabilidad de que se necesite cambiarle el aceite y el filtro es 0.14.

- Si se le tiene que cambiar el aceite, ¿cuál es la probabilidad de que también se necesite cambiarle el filtro?
- Si se le tiene que cambiar el filtro de aceite, ¿cuál es la probabilidad de que también se le tenga que cambiar el aceite?

2.81 La probabilidad de que un hombre casado vea cierto programa de televisión es 0.4 y la probabilidad de que lo vea una mujer casada es 0.5. La probabilidad

de que un hombre vea el programa, dado que su esposa lo ve, es 0.7. Calcule la probabilidad de que

- una pareja casada vea el programa;
- una esposa vea el programa dado que su esposo lo ve;
- al menos uno de los miembros de la pareja casada vea el programa.

2.82 Para parejas casadas que viven en cierto suburbio, la probabilidad de que el esposo vote en un referéndum es 0.21, la probabilidad de que vote la esposa es 0.28 y la probabilidad de que ambos voten es 0.15. ¿Cuál es la probabilidad de que...

- al menos uno de los miembros de la pareja casada vote?
- una esposa vote, dado que su esposo vota?
- un esposo vote, dado que su esposa no vota?

2.83 La probabilidad de que un vehículo que entra a las Cavernas Luray tenga matrícula de Canadá es 0.12, la probabilidad de que sea una casa rodante es 0.28 y la probabilidad de que sea una casa rodante con matrícula de Canadá es 0.09. ¿Cuál es la probabilidad de que...

- una casa rodante que entra a las Cavernas Luray tenga matrícula de Canadá?
- un vehículo con matrícula de Canadá que entra a las Cavernas Luray sea una casa rodante?
- un vehículo que entra a las Cavernas Luray no tenga matrícula de Canadá o no sea una casa rodante?

2.84 La probabilidad de que el jefe de familia esté en casa cuando llame el representante de marketing de una empresa es 0.4. Dado que el jefe de familia está en casa, la probabilidad de que la empresa le venda un producto es 0.3. Encuentre la probabilidad de que el jefe de familia esté en casa y compre productos de la empresa.

2.85 La probabilidad de que un doctor diagnostique de manera correcta una enfermedad específica es 0.7. Dado que el doctor hace un diagnóstico incorrecto, la probabilidad de que el paciente entable una demanda legal es 0.9. ¿Cuál es la probabilidad de que el doctor haga un diagnóstico incorrecto y el paciente lo demande?

2.86 En 1970, 11% de los estadounidenses completaron cuatro años de universidad; de ese porcentaje 43 % eran mujeres. En 1990, 22% de los estadounidenses completaron cuatro años de universidad, un porcentaje del cual 53 % fueron mujeres. (*Time*, 19 de enero de 1996).

- Dado que una persona completó cuatro años de universidad en 1970, ¿cuál es la probabilidad de que esa persona sea mujer?

- b) ¿Cuál es la probabilidad de que una mujer haya terminado cuatro años de universidad en 1990?
- c) ¿Cuál es la probabilidad de que en 1990 un hombre no haya terminado la universidad?

2.87 Un agente de bienes raíces tiene 8 llaves maestras para abrir varias casas nuevas. Sólo 1 llave maestra abrirá cualquiera de las casas. Si 40% de estas casas por lo general se dejan abiertas, ¿cuál es la probabilidad de que el agente de bienes raíces pueda entrar en una casa específica, si selecciona 3 llaves maestras al azar antes de salir de la oficina?

2.88 Antes de la distribución de cierto software estadístico se prueba la precisión de cada cuarto disco compacto (CD). El proceso de prueba consiste en correr cuatro programas independientes y verificar los resultados. La tasa de falla para los 4 programas de prueba son 0.01, 0.03, 0.02 y 0.01, respectivamente.

- a) ¿Cuál es la probabilidad de que uno de los CD que se pruebe no pase la prueba?
- b) Dado que se prueba un CD, ¿cuál es la probabilidad de que falle el programa 2 o 3?
- c) En una muestra de 100, ¿cuántos CD esperarías que se rechazaran?
- d) Dado que un CD está defectuoso, ¿cuál es la probabilidad de que se pruebe?

2.89 Una ciudad tiene dos carros de bomberos que operan de forma independiente. La probabilidad de que un carro específico esté disponible cuando se le necesite es 0.96.

- a) ¿Cuál es la probabilidad de que ninguno esté disponible cuando se necesite?
- b) ¿Cuál es la probabilidad de que un carro de bomberos esté disponible cuando se le necesite?

2.90 La contaminación de los ríos de Estados Unidos ha sido un problema por muchos años. Considere los siguientes eventos:

- A: el río está contaminado.
 B: al probar una muestra de agua se detecta contaminación.
 C: se permite pescar.

Suponga que $P(A) = 0.3$, $P(B|A) = 0.75$, $P(B|A^c) = 0.20$, $P(C|A \cap B) = 0.20$, $P(C|A^c \cap B) = 0.15$, $P(C|A \cap B^c) = 0.80$ y $P(C|A^c \cap B^c) = 0.90$.

- a) Calcule $P(A \cap B \cap C)$.
 b) Calcule $P(B^c \cap C)$.
 c) Calcule $P(C)$.
 d) Calcule la probabilidad de que el río esté contaminado, dado que está permitido pescar y que la muestra probada no detectó contaminación.

2.91 Encuentre la posibilidad de seleccionar aleatoriamente 4 litros de leche en buenas condiciones sucesivamente de un refrigerador que contiene 20 litros, de los cuales 5 están echados a perder, utilizando

- a) la primera fórmula del teorema 2.12 de la página 68;
 b) las fórmulas del teorema 2.6 y la regla 2.3 de las páginas 50 y 54, respectivamente.

2.92 Imagine el diagrama de un sistema eléctrico como el que se muestra en la figura 2.10. ¿Cuál es la probabilidad de que el sistema funcione? Suponga que los componentes fallan de forma independiente.

2.93 En la figura 2.11 se muestra un sistema de circuitos. Suponga que los componentes fallan de manera independiente.

- a) ¿Cuál es la probabilidad de que el sistema completo funcione?
 b) Dado que el sistema funciona, ¿cuál es la probabilidad de que el componente A no funcione?

2.94 En la situación del ejercicio 2.93 se sabe que el sistema no funciona. ¿Cuál es la probabilidad de que el componente A tampoco funcione?

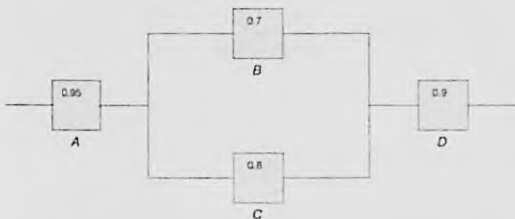


Figura 2.10: Diagrama para el ejercicio 2.92.

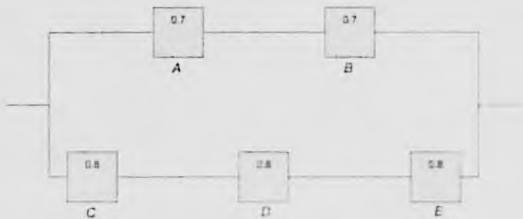


Figura 2.11: Diagrama para el ejercicio 2.93.

2.7 Regla de Bayes

La estadística bayesiana es un conjunto de herramientas que se utiliza en un tipo especial de inferencia estadística que se aplica en el análisis de datos experimentales en muchas situaciones prácticas de ciencia e ingeniería. La regla de Bayes es una de las normas más importantes de la teoría de probabilidad, ya que es el fundamento de la inferencia bayesiana, la cual se analizará en el capítulo 18.

Probabilidad total

Regresemos al ejemplo de la sección 2.6, en el que se selecciona un individuo al azar de entre los adultos de una pequeña ciudad para que viaje por el país promoviendo las ventajas de establecer industrias nuevas en la ciudad. Suponga que ahora se nos da la información adicional de que 36 de los empleados y 12 de los desempleados son miembros del Club Rotario. Deseamos encontrar la probabilidad del evento A de que el individuo seleccionado sea miembro del Club Rotario. Podemos remitirnos a la figura 2.12 y escribir A como la unión de los dos eventos mutuamente excluyentes $E \cap A$ y $E' \cap A$. Por lo tanto, $A = (E \cap A) \cup (E' \cap A)$, y mediante el corolario 2.1 del teorema 2.7 y luego mediante el teorema 2.10, podemos escribir

$$\begin{aligned} P(A) &= P[(E \cap A) \cup (E' \cap A)] = P(E \cap A) + P(E' \cap A) \\ &= P(E)P(A|E) + P(E')P(A|E'). \end{aligned}$$

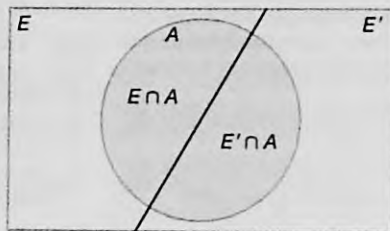


Figura 2.12: Diagrama de Venn para los eventos A , E y E' .

Los datos de la sección 2.6, junto con los datos adicionales antes dados para el conjunto A , nos permiten calcular

$$P(E) = \frac{600}{900} = \frac{2}{3}, \quad P(A|E) = \frac{36}{600} = \frac{3}{50},$$

y

$$P(E') = \frac{1}{3}, \quad P(A|E') = \frac{12}{300} = \frac{1}{25}.$$

Si mostramos estas probabilidades mediante el diagrama de árbol de la figura 2.13, donde la primera rama da la probabilidad $P(E)P(A|E)$ y la segunda rama da la probabilidad

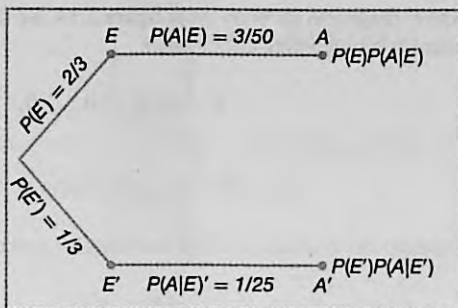


Figura 2.13: Diagrama de árbol para los datos de la página 63 con información adicional de la página 72.

la probabilidad $P(E')P(A|E')$, deducimos que

$$P(A) = \left(\frac{2}{3}\right) \left(\frac{3}{50}\right) + \left(\frac{1}{3}\right) \left(\frac{1}{25}\right) = \frac{4}{75}.$$

Una generalización del ejemplo anterior para el caso en donde el espacio muestral se parte en k subconjuntos se cubre mediante el siguiente teorema, que algunas veces se denomina **teorema de probabilidad total** o **regla de eliminación**.

Teorema 2.13: Si los eventos B_1, B_2, \dots, B_k constituyen una partición del espacio muestral S , tal que $P(B_i) \neq 0$ para $i = 1, 2, \dots, k$, entonces, para cualquier evento A de S ,

$$P(A) = \sum_{i=1}^k P(B_i \cap A) = \sum_{i=1}^k P(B_i)P(A|B_i).$$

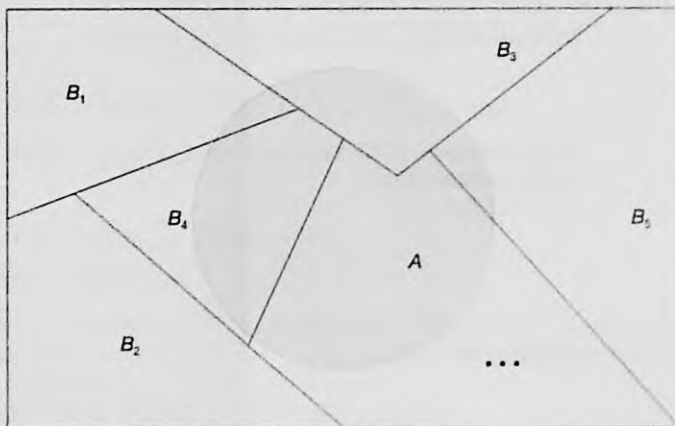


Figura 2.14: Partición del espacio muestral s .

Prueba: Considere el diagrama de Venn de la figura 2.14. Se observa que el evento A es la unión de los eventos mutuamente excluyentes

$$B_1 \cap A, B_2 \cap A, \dots, B_k \cap A;$$

es decir,

$$A = (B_1 \cap A) \cup (B_2 \cap A) \cup \dots \cup (B_k \cap A)$$

Por medio del corolario 2.2 del teorema 2.7 y el teorema 2.10 obtenemos

$$\begin{aligned} P(A) &= P[(B_1 \cap A) \cup (B_2 \cap A) \cup \dots \cup (B_k \cap A)] \\ &= P(B_1 \cap A) + P(B_2 \cap A) + \dots + P(B_k \cap A) \\ &= \sum_{i=1}^k P(B_i \cap A) \\ &= \sum_{i=1}^k P(B_i)P(A|B_i). \end{aligned}$$

Ejemplo 2.41: Tres máquinas de cierta planta de ensamble, B_1 , B_2 y B_3 , montan 30%, 45% y 25% de los productos, respectivamente. Se sabe por experiencia que 2%, 3% y 2% de los productos ensamblados por cada máquina, respectivamente, tienen defectos. Ahora bien, suponga que se selecciona de forma aleatoria un producto terminado. ¿Cuál es la probabilidad de que esté defectuoso?

Solución: Considere los siguientes eventos:

A : el producto está defectuoso,

B_1 : el producto fue ensamblado con la máquina B_1 ,

B_2 : el producto fue ensamblado con la máquina B_2 ,

B_3 : el producto fue ensamblado con la máquina B_3 .

Podemos aplicar la regla de eliminación y escribir

$$P(A) = P(B_1)P(A|B_1) + P(B_2)P(A|B_2) + P(B_3)P(A|B_3).$$

Si nos remitimos al diagrama de árbol de la figura 2.15 encontramos que las tres ramas dan las probabilidades

$$P(B_1)P(A|B_1) = (0.3)(0.02) = 0.006,$$

$$P(B_2)P(A|B_2) = (0.45)(0.03) = 0.0135,$$

$$P(B_3)P(A|B_3) = (0.25)(0.02) = 0.005,$$

en consecuencia,

$$P(A) = 0.006 + 0.0135 + 0.005 = 0.0245.$$

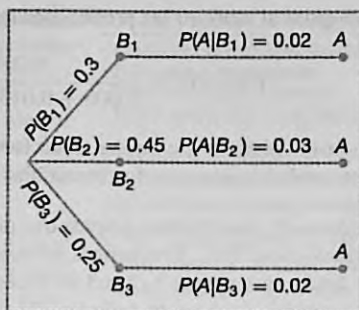


Figura 2.15: Diagrama de árbol para el ejemplo 2.41.

Regla de Bayes

Suponga que en lugar de calcular $P(A)$ mediante la regla de eliminación en el ejemplo 2.41, consideramos el problema de obtener la probabilidad condicional $P(B_i|A)$. En otras palabras, suponga que se selecciona un producto de forma aleatoria y que éste resulta defectuoso. ¿Cuál es la probabilidad de que este producto haya sido ensamblado con la máquina B_i ? Las preguntas de este tipo se pueden contestar usando el siguiente teorema, denominado **regla de Bayes**:

Teorema 2.14: (Regla de Bayes) Si los eventos B_1, B_2, \dots, B_k constituyen una partición del espacio muestral S , donde $P(B_i) \neq 0$ para $i = 1, 2, \dots, k$, entonces, para cualquier evento A en S , tal que $P(A) \neq 0$,

$$P(B_r|A) = \frac{P(B_r \cap A)}{\sum_{i=1}^k P(B_i \cap A)} = \frac{P(B_r)P(A|B_r)}{\sum_{i=1}^k P(B_i)P(A|B_i)} \quad \text{para } r = 1, 2, \dots, k.$$

Prueba: Mediante la definición de probabilidad condicional,

$$P(B_r|A) = \frac{P(B_r \cap A)}{P(A)},$$

y después usando el teorema 2.13 en el denominador, tenemos

$$P(B_r|A) = \frac{P(B_r \cap A)}{\sum_{i=1}^k P(B_i \cap A)} = \frac{P(B_r)P(A|B_r)}{\sum_{i=1}^k P(B_i)P(A|B_i)},$$

que completa la demostración. ▮

Ejemplo 2.42: Con referencia al ejemplo 2.41, si se elige al azar un producto y se encuentra que está defectuoso, ¿cuál es la probabilidad de que haya sido ensamblado con la máquina B_3 ?

Solución: Podemos utilizar la regla de Bayes para escribir

$$P(B_3|A) = \frac{P(B_3)P(A|B_3)}{P(B_1)P(A|B_1) + P(B_2)P(A|B_2) + P(B_3)P(A|B_3)},$$

y después al sustituir las probabilidades calculadas en el ejemplo 2.41, tenemos

$$P(B_3|A) = \frac{0.005}{0.006 + 0.0135 + 0.005} = \frac{0.005}{0.0245} = \frac{10}{49}$$

En vista del hecho de que se seleccionó un producto defectuoso, este resultado sugiere que probablemente no fue ensamblado con la máquina B_3 . \blacksquare

Ejemplo 2.43: Una empresa de manufactura emplea tres planos analíticos para el diseño y desarrollo de un producto específico. Por razones de costos los tres se utilizan en momentos diferentes. De hecho, los planos 1, 2 y 3 se utilizan para 30%, 20% y 50% de los productos, respectivamente. La tasa de defectos difiere en los tres procedimientos de la siguiente manera,

$$P(D|P_1) = 0.01, \quad P(D|P_2) = 0.03, \quad P(D|P_3) = 0.02,$$

en donde $P(D|P_j)$ es la probabilidad de que un producto esté defectuoso, dado el plano j . Si se observa un producto al azar y se descubre que está defectuoso, ¿cuál de los planos tiene más probabilidades de haberse utilizado y, por lo tanto, de ser el responsable?

Solución: A partir del planteamiento del problema

$$P(P_1) = 0.30, \quad P(P_2) = 0.20 \quad \text{y} \quad P(P_3) = 0.50,$$

debemos calcular $P(P_j|D)$ para $j = 1, 2, 3$. La regla de Bayes (teorema 2.14) muestra

$$\begin{aligned} P(P_1|D) &= \frac{P(P_1)P(D|P_1)}{P(P_1)P(D|P_1) + P(P_2)P(D|P_2) + P(P_3)P(D|P_3)} \\ &= \frac{(0.30)(0.01)}{(0.3)(0.01) + (0.20)(0.03) + (0.50)(0.02)} = \frac{0.003}{0.019} = 0.158. \end{aligned}$$

De igual manera,

$$P(P_2|D) = \frac{(0.03)(0.20)}{0.019} = 0.316 \quad \text{y} \quad P(P_3|D) = \frac{(0.02)(0.50)}{0.019} = 0.526.$$

La probabilidad condicional de un defecto, dado el plano 3, es la mayor de las tres; por consiguiente, un defecto en un producto elegido al azar tiene más probabilidad de ser el resultado de haber usado el plano 3. \blacksquare

La regla de Bayes, un método estadístico llamado método bayesiano, ha adquirido muchas aplicaciones. En el capítulo 18 estudiaremos una introducción al método bayesiano.

Ejercicios

2.95 En cierta región del país se sabe por experiencia que la probabilidad de seleccionar un adulto mayor de 40 años de edad con cáncer es 0.05. Si la probabilidad de que un doctor diagnostique de forma correcta que una persona con cáncer tiene la enfermedad es 0.78, y la probabilidad de que diagnostique de forma incorrecta que una persona sin cáncer tiene la enfermedad es 0.06, ¿cuál es la probabilidad de que a un adulto mayor de 40 años se le diagnostique cáncer?

2.96 La policía planea hacer respetar los límites de velocidad usando un sistema de radar en 4 diferentes puntos a las orillas de la ciudad. Las trampas de radar en cada uno de los sitios L_1 , L_2 , L_3 y L_4 operarán 40%, 30%, 20% y 30% del tiempo. Si una persona que excede el límite de velocidad cuando va a su trabajo tiene probabilidades de 0.2, 0.1, 0.5 y 0.2, respectivamente, de pasar por esos lugares, ¿cuál es la probabilidad de que reciba una multa por conducir con exceso de velocidad?

2.97 Remítase al ejercicio 2.95. ¿Cuál es la probabilidad de que una persona a la que se le diagnostica cáncer realmente tenga la enfermedad?

2.98 Si en el ejercicio 2.96 la persona es multada por conducir con exceso de velocidad en su camino al trabajo, ¿cuál es la probabilidad de que pase por el sistema de radar que se ubica en L_2 ?

2.99 Suponga que los cuatro inspectores de una fábrica de película colocan la fecha de caducidad en cada paquete de película al final de la línea de montaje. John, quien coloca la fecha de caducidad en 20% de los paquetes, no logra ponerla en uno de cada 200 paquetes; Tom, quien la coloca en 60% de los paquetes, no logra ponerla en uno de cada 100 paquetes; Jeff, quien la coloca en 15% de los paquetes, no lo hace una vez en cada 90 paquetes; y Pat, que fecha 5% de los paquetes, falla en uno de cada 200 paquetes. Si un cliente se queja de que su paquete de película no muestra la fecha de caducidad, ¿cuál es la probabilidad de que haya sido inspeccionado por John?

2.100 Una empresa telefónica regional opera tres estaciones de retransmisión idénticas en diferentes sitios. A continuación se muestra el número de desperfectos en cada estación reportados durante un año y las causas de éstos.

	Estación A	B	C
Problemas con el suministro de electricidad	2	1	1
Falla de la computadora	4	3	2
Fallas del equipo eléctrico	5	4	2
Fallas ocasionadas por otros errores humanos	7	5	5

Suponga que se reporta una falla y que se descubre que fue ocasionada por otros errores humanos. ¿Cuál es la probabilidad de que provenga de la estación C?

2.101 Una cadena de tiendas de pintura produce y vende pintura de látex y semiesmaltada. De acuerdo con las ventas a largo plazo, la probabilidad de que un cliente compre pintura de látex es 0.75. De los que compran pintura de látex, 60 % también compra rodillos. Sin embargo, sólo 30 % de los que compran pintura semiesmaltada compra rodillos. Un comprador que se selecciona al azar adquiere un rodillo y una lata de pintura. ¿Cuál es la probabilidad de que sea pintura de látex?

2.102 Denote como A , B y C a los eventos de que un gran premio se encuentra detrás de las puertas A , B y C , respectivamente. Suponga que elige al azar una puerta, por ejemplo la A . El presentador del juego abre una puerta, por ejemplo la B , y muestra que no hay un premio detrás de ella. Ahora, el presentador le da la opción de conservar la puerta que eligió (A) o de cambiarla por la puerta que queda (C). Utilice la probabilidad para explicar si debe o no hacer el cambio.

Ejercicios de repaso

2.103 Un suero de la verdad tiene la propiedad de que 90% de los sospechosos culpables se juzgan de forma adecuada, mientras que, por supuesto, 10% de los sospechosos culpables erróneamente se consideran inocentes. Por otro lado, a los sospechosos inocentes se les juzga de manera errónea 1% de las veces. Si se aplica el suero a un sospechoso, que se selecciona de un grupo de sospechosos en el cual sólo 5% ha cometido un delito, y éste indica que es culpable, ¿cuál es la probabilidad de que sea inocente?

2.104 Un alergólogo afirma que 50% de los pacientes que examina son alérgicos a algún tipo de hierba. ¿Cuál es la probabilidad de que...

- exactamente 3 de sus 4 pacientes siguientes sean alérgicos a hierbas?
- ninguno de sus 4 pacientes siguientes sea alérgico a hierbas?

2.105 Mediante la comparación de las regiones apropiadas en un diagrama de Venn, verifique que

- $(A \cap B) \cup (A \cap B') = A$;
- $A' \cap (B' \cup C) = (A' \cap B') \cup (A' \cap C)$.

2.106 Las probabilidades de que una estación de servicio bombee gasolina en 0, 1, 2, 3, 4, 5 o más automóviles durante cierto periodo de 30 minutos son, respectivamente, 0.03, 0.18, 0.24, 0.28, 0.10 y 0.17. Calcule la probabilidad de que en este periodo de 30 minutos

- más de 2 automóviles reciban gasolina;
- a lo sumo 4 automóviles reciban gasolina;
- 4 o más automóviles reciban gasolina.

2.107 ¿Cuántas manos de *bridge* que contengan 4 espadas, 6 diamantes, 1 trébol y 2 corazones son posibles?

2.108 Si la probabilidad de que una persona cometa un error en su declaración de impuestos sobre la renta es 0.1, calcule la probabilidad de que

- cada una de cuatro personas no relacionadas cometa un error;
- el señor Jones y la señora Clark cometan un error, y el señor Roberts y la señora Williams no cometan errores.

- 2.109** Una empresa industrial grande usa tres moteles locales para ofrecer hospedaje nocturno a sus clientes. Se sabe por experiencia que a 20% de los clientes se le asigna habitaciones en el Ramada Inn, a 50% en el Sheraton y a 30% en el Lakeview Motor Lodge. Si hay una falla en la plomería en 5% de las habitaciones del Ramada Inn, en 4% de las habitaciones del Sheraton y en 8% de las habitaciones del Lakeview Motor Lodge, ¿cuál es la probabilidad de que...
- a) a un cliente se le asigne una habitación en la que falle la plomería?
 - b) a una persona que ocupa una habitación en la que falla la plomería se le haya hospedado en el Lakeview Motor Lodge?
- 2.110** La probabilidad de que un paciente se recupere de una delicada operación de corazón es 0.8. ¿Cuál es la probabilidad de que...
- a) exactamente 2 de los siguientes 3 pacientes a los que se somete a esta operación sobrevivan?
 - b) los siguientes 3 pacientes que tengan esta operación sobrevivan?
- 2.111** Se sabe que $\frac{2}{3}$ de los reclusos en cierta prisión federal son menores de 25 años de edad. También se sabe que $\frac{3}{5}$ de los reos son hombres y que $\frac{5}{8}$ son mujeres de 25 años de edad o mayores. ¿Cuál es la probabilidad de que un prisionero seleccionado al azar de esta prisión sea mujer y tenga al menos 25 años de edad?
- 2.112** Si se tienen 4 manzanas rojas, 5 verdes y 6 amarillas, ¿cuántas selecciones de 9 manzanas se pueden hacer si se deben seleccionar 3 de cada color?
- 2.113** De una caja que contiene 6 bolas negras y 4 verdes se extraen 3 bolas sucesivamente y cada bola se reemplaza en la caja antes de extraer la siguiente. ¿Cuál es la probabilidad de que...
- a) las 3 sean del mismo color?
 - b) cada color esté representado?
- 2.114** Un cargamento de 12 televisores contiene tres defectuosos. ¿De cuántas formas puede un hotel comprar 5 de estos aparatos y recibir al menos 2 defectuosos?
- 2.115** Cierta organización federal emplea a tres empresas consultoras (A, B y C) con probabilidades de 0.40, 0.35 y 0.25, respectivamente. Se sabe por experiencia que las probabilidades de que las empresas rebasen los costos son 0.05, 0.03 y 0.15, respectivamente. Suponga que el organismo experimenta un exceso en los costos.
- a) ¿Cuál es la probabilidad de que la empresa consultora implicada sea la C?
 - b) ¿Cuál es la probabilidad de que sea la A?
- 2.116** Un fabricante estudia los efectos de la temperatura de cocción, el tiempo de cocción y el tipo de aceite para la cocción al elaborar papas fritas. Se utilizan 3 diferentes temperaturas, 4 diferentes tiempos de cocción y 3 diferentes aceites.
- a) ¿Cuál es el número total de combinaciones a estudiar?
 - b) ¿Cuántas combinaciones se utilizarán para cada tipo de aceite?
 - c) Analice por qué las permutaciones no intervienen en este ejercicio.
- 2.117** Considere la situación del ejercicio 2.116 y suponga que el fabricante puede probar sólo dos combinaciones en un día.
- a) ¿Cuál es la probabilidad de que elija cualquier conjunto dado de 2 corridas?
 - b) ¿Cuál es la probabilidad de que utilice la temperatura más alta en cualquiera de estas 2 combinaciones?
- 2.118** Se sabe que existe una probabilidad de 0.07 de que las mujeres de más de 60 años desarrollen cierta forma de cáncer. Se dispone de una prueba de sangre que, aunque no es infalible, permite detectar la enfermedad. De hecho, se sabe que 10 % de las veces la prueba da un falso negativo (es decir, la prueba da un resultado negativo de manera incorrecta) y 5 % de las veces la prueba da un falso positivo (es decir, la prueba da un resultado positivo de manera incorrecta). Si una mujer de más de 60 años se somete a la prueba y recibe un resultado favorable (es decir, negativo), ¿qué probabilidad hay de que tenga la enfermedad?
- 2.119** Un fabricante de cierto tipo de componente electrónico abastece a los proveedores en lotes de 20. Suponga que 60% de todos los lotes no contiene componentes defectuosos, que 30% contiene un componente defectuoso y que 10% contiene dos componentes defectuosos. Si se elige un lote del que se extraen aleatoriamente dos componentes, los cuales se prueban y ninguno resulta defectuoso,
- a) ¿Cuál es la probabilidad de que haya cero componentes defectuosos en el lote?
 - b) ¿Cuál es la probabilidad de que haya un componente defectuoso en el lote?
 - c) ¿Cuál es la probabilidad de que haya dos componentes defectuosos en el lote?
- 2.120** Existe una extraña enfermedad que sólo afecta a uno de cada 500 individuos. Se dispone de una prueba para detectarla, pero, por supuesto, ésta no es infalible. Un resultado correcto positivo (un paciente que realmente tiene la enfermedad) ocurre 95% de las veces; en tanto que un resultado falso positivo (un paciente que no tiene la enfermedad) ocurre 1% de las veces. Si un individuo elegido al azar se somete a prueba y se obtiene un resultado positivo, ¿cuál es la probabilidad de que realmente tenga la enfermedad?
- 2.121** Una empresa constructora emplea a dos ingenieros de ventas. El ingeniero 1 hace el trabajo de estimar costos en 70% de las cotizaciones solicitadas a la empresa. El ingeniero 2 hace lo mismo en 30% de las

cotizaciones. Se sabe que la tasa de error para el ingeniero 1 es tal que la probabilidad de encontrar un error en su trabajo es 0.02; mientras que la probabilidad de encontrar un error en el trabajo del ingeniero 2 es 0.04. Suponga que al revisar una solicitud de cotización se encuentra un error grave en la estimación de los costos. ¿Qué ingeniero supondría usted que hizo los cálculos? Explique su respuesta y muestre todo el desarrollo.

2.122 En el campo del control de calidad a menudo se usa la ciencia estadística para determinar si un proceso está "fuera de control". Suponga que el proceso, de hecho, está fuera de control y que 20 por ciento de los artículos producidos tiene defecto.

- Si tres artículos salen en serie de la línea de producción, ¿cuál es la probabilidad de que los tres estén defectuosos?
- Si salen cuatro artículos en serie, ¿cuál es la probabilidad de que tres estén defectuosos?

2.123 En una planta industrial se está realizando un estudio para determinar la rapidez con la que los trabajadores lesionados regresan a sus labores después del percance. Los registros demuestran que 10% de los trabajadores lesionados son llevados al hospital para su tratamiento y que 15% regresan a su trabajo al día siguiente. Además, los estudios demuestran que 2% son llevados al hospital y regresan al trabajo al día siguiente. Si un trabajador se lesiona, ¿cuál es la probabilidad de que sea llevado al hospital, de que regrese al trabajo al día siguiente, o de ambas cosas?

2.124 Una empresa acostumbra capacitar operadores que realizan ciertas actividades en la línea de producción. Se sabe que los operadores que asisten al curso de capacitación son capaces de cumplir sus cuotas de producción 90% de las veces. Los nuevos operarios que no toman el curso de capacitación sólo cumplen con sus cuotas 65% de las veces. Cincuenta por ciento de los nuevos operadores asisten al curso. Dado que un nuevo operador cumple con su cuota de producción, ¿cuál es la probabilidad de que haya asistido al curso?

2.125 Una encuesta aplicada a quienes usan un software estadístico específico indica que 10% no quedó satisfecho. La mitad de quienes no quedaron satisfechos le compraron el sistema al vendedor A. También se sabe que 20% de los encuestados se lo compraron al

vendedor A. Dado que el proveedor del paquete de software fue el vendedor A, ¿cuál es la probabilidad de que un usuario específico haya quedado insatisfecho?

2.126 Durante las crisis económicas se despide a obreros y a menudo se les reemplaza con máquinas. Se revisa la historia de 100 trabajadores cuya pérdida del empleo se atribuye a los avances tecnológicos. Para cada uno de ellos se determinó si obtuvieron un empleo alternativo dentro de la misma empresa, si encontraron un empleo en la misma área de otra empresa, si encontraron trabajo en una nueva área o si llevan desempleados más de un año. Además, se registró la situación sindical de cada trabajador. La siguiente tabla resume los resultados.

	No Sindicalizado sindicalizado	
Está en la misma empresa	40	15
Está en otra empresa (misma área)	13	10
Está en una nueva área	4	11
Está desempleado	2	5

- Si un trabajador seleccionado encontró empleo en la misma área de una nueva empresa, ¿cuál es la probabilidad de que sea miembro de un sindicato?
- Si el trabajador es miembro de un sindicato, ¿cuál es la probabilidad de que esté desempleado desde hace un año?

2.127 Hay 50% de probabilidad de que la reina tenga el gen de la hemofilia. Si es portadora, entonces cada uno de los príncipes tiene 50% de probabilidad independiente de tener hemofilia. Si la reina no es portadora, el príncipe no tendrá la enfermedad. Suponga que la reina tuvo tres príncipes que no padecen la enfermedad, ¿cuál es la probabilidad de que la reina sea portadora del gen?

2.128 Proyecto de equipo: Entregue a cada estudiante una bolsa de chocolates M&M y forme equipos de 5 o 6 estudiantes. Calcule la distribución de frecuencia relativa del color de los M&M para cada equipo.

- ¿Cuál es su probabilidad estimada de elegir un chocolate amarillo al azar? ¿Y uno rojo?
- Ahora haga el mismo cálculo para todo el grupo. ¿Cambiaron las estimaciones?
- ¿Cree que en un lote procesado existe el mismo número de chocolates de cada color? Comente al respecto.

2.8 Posibles riesgos y errores conceptuales; relación con el material de otros capítulos

Este capítulo incluye las definiciones, reglas y teoremas fundamentales que convierten a la probabilidad en una herramienta importante para la evaluación de sistemas científicos y de ingeniería. A menudo estas evaluaciones toman la forma de cálculos de probabili-

dad, como se ilustra en los ejemplos y en los ejercicios. Conceptos como independencia, probabilidad condicional, regla de Bayes y otros suelen ser muy adecuados para resolver problemas prácticos en los que se busca obtener un valor de probabilidad. Abundan las ilustraciones en los ejercicios. Vea, por ejemplo, los ejercicios 2.100 y 2.101. En éstos y en muchos otros ejercicios se realiza una evaluación juiciosa de un sistema científico, a partir de un cálculo de probabilidad, utilizando las reglas y las definiciones que se estudian en el capítulo.

Ahora bien, ¿qué relación existe entre el material de este capítulo y el material de otros capítulos? La mejor forma de responder esta pregunta es dando un vistazo al capítulo 3, ya que en éste también se abordan problemas en los que es importante el cálculo de probabilidades. Ahí se ilustra cómo el desempeño de un sistema depende del valor de una o más probabilidades. De nuevo, la probabilidad condicional y la independencia desempeñan un papel. Sin embargo, surgen nuevos conceptos que permiten tener una mayor estructura basada en el concepto de una variable aleatoria y su distribución de probabilidad. Recuerde que el concepto de las distribuciones de frecuencias se abordó brevemente en el capítulo 1. La distribución de probabilidad muestra, en forma gráfica o en una ecuación, toda la información necesaria para describir una estructura de probabilidad. Por ejemplo, en el ejercicio de repaso 2.122 la variable aleatoria de interés es el número de artículos defectuosos, una medición discreta. Por consiguiente, la distribución de probabilidad revelaría la estructura de probabilidad para el número de artículos defectuosos extraídos del número elegido del proceso. Cuando el lector avance al capítulo 3 y los siguientes, será evidente para él que se requieren suposiciones para determinar y, por lo tanto, utilizar las distribuciones de probabilidad en la resolución de problemas científicos.

CAPÍTULO 3

Variables aleatorias y distribuciones de probabilidad

3.1 Concepto de variable aleatoria

La estadística realiza inferencias acerca de las poblaciones y sus características. Se llevan a cabo experimentos cuyos resultados se encuentran sujetos al azar. La prueba de un número de componentes electrónicos es un ejemplo de **experimento estadístico**, un concepto que se utiliza para describir cualquier proceso mediante el cual se generan varias observaciones al azar. A menudo es importante asignar una descripción numérica al resultado. Por ejemplo, cuando se prueban tres componentes electrónicos, el espacio muestral que ofrece una descripción detallada de cada posible resultado se escribe como

$$S = \{NNN, NND, NDN, DNN, NDD, DND, DDN, DDD\},$$

donde N denota “no defectuoso”, y D , “defectuoso”. Es evidente que nos interesa el número de componentes defectuosos que se presenten. De esta forma, a cada punto en el espacio muestral se le *asignará un valor numérico* de 0, 1, 2 o 3. Estos valores son, por supuesto, cantidades aleatorias *determinadas por el resultado del experimento*. Se pueden ver como valores que toma la *variable aleatoria* X , es decir, el número de artículos defectuosos cuando se prueban tres componentes electrónicos.

Definición 3.1: Una **variable aleatoria** es una función que asocia un número real con cada elemento del espacio muestral.

Utilizaremos una letra mayúscula, digamos X , para denotar una variable aleatoria, y su correspondiente letra minúscula, x en este caso, para uno de sus valores. En el ejemplo de la prueba de componentes electrónicos observamos que la variable aleatoria X toma el valor 2 para todos los elementos en el subconjunto

$$E = \{DDN, DND, NDD\}$$

del espacio muestral S . Esto es, cada valor posible de X representa un evento que es un subconjunto del espacio muestral para el experimento dado.

Ejemplo 3.1: De una urna que contiene 4 bolas rojas y 3 negras se sacan 2 bolas de manera sucesiva, sin reemplazo. Los posibles resultados y los valores y de la variable aleatoria Y , donde Y es el número de bolas rojas, son

Espacio muestral	y
RR	2
RN	1
NR	1
NN	0

Ejemplo 3.2: El empleado de un almacén regresa tres cascos de seguridad al azar a tres obreros de un taller siderúrgico que ya los habían probado. Si Smith, Jones y Brown, en ese orden, reciben uno de los tres cascos, liste los puntos muestrales para los posibles órdenes en que el empleado del almacén regresa los cascos, después calcule el valor m de la variable aleatoria M que representa el número de emparejamientos correctos.

Solución: Si S , J y B representan, respectivamente, los cascos que recibieron Smith, Jones y Brown, entonces los posibles arreglos en los cuales se pueden regresar los cascos y el número de emparejamientos correctos son

Espacio muestral	m
SJB	3
SBJ	1
BJS	1
JSB	1
JBS	0
BSJ	0

En cada uno de los dos ejemplos anteriores, el espacio muestral contiene un número finito de elementos. Por el contrario, cuando lanzamos un dado hasta que salga un 5, obtenemos un espacio muestral con una secuencia de elementos interminable,

$$S = \{F, NF, NNF, NNNF, \dots\},$$

donde F y N representan, respectivamente, la ocurrencia y la no ocurrencia de un 5. Sin embargo, incluso en este experimento el número de elementos se puede igualar a la cantidad total de números enteros, de manera que hay un primer elemento, un segundo, un tercero y así sucesivamente, por lo que se pueden contar.

Hay casos en que la variable aleatoria es categórica por naturaleza en los cuales se utilizan las llamadas variables *ficticias*. Un buen ejemplo de ello es el caso en que la variable aleatoria es binaria por naturaleza, como se indica a continuación.

Ejemplo 3.3: Considere la condición en que salen componentes de la línea de ensamble y se les clasifica como defectuosos o no defectuosos. Defina la variable aleatoria X mediante

$$X = \begin{cases} 1, & \text{si el componente está defectuoso,} \\ 0, & \text{si el componente no está defectuoso.} \end{cases}$$

Evidentemente la asignación de 1 o 0 es arbitraria, aunque bastante conveniente, lo cual quedará más claro en capítulos posteriores. La variable aleatoria en la que se eligen 0 y 1 para describir los dos posibles valores se denomina **variable aleatoria de Bernoulli**. ■

En los siguientes ejemplos veremos más casos de variables aleatorias.

Ejemplo 3.4: Los estadísticos utilizan **planes de muestreo** para aceptar o rechazar lotes de materiales. Suponga que uno de los planes de muestreo implica obtener una muestra independiente de 10 artículos de un lote de 100, en el que 12 están defectuosos.

Si X representa a la variable aleatoria, definida como el número de artículos que están defectuosos en la muestra de 10, la variable aleatoria toma los valores 0, 1, 2, ..., 9, 10. ■

Ejemplo 3.5: Suponga que un plan de muestreo implica obtener una muestra de artículos de un proceso hasta que se encuentre uno defectuoso. La evaluación del proceso dependerá de cuántos artículos consecutivos se observen. En este caso, sea X una variable aleatoria que se define como el número de artículos observados antes de que salga uno defectuoso. Si N representa un artículo no defectuoso y D uno defectuoso, los espacios muestrales son $S = \{D\}$ dado que $X = 1$, $S = \{ND\}$ dado que $X = 2$, $S = \{NND\}$ dado que $X = 3$, y así sucesivamente. ■

Ejemplo 3.6: Existe interés por la proporción de personas que responden a cierta encuesta enviada por correo. Sea X tal proporción. X es una variable aleatoria que toma todos los valores de x para los cuales $0 \leq x \leq 1$. ■

Ejemplo 3.7: Sea X la variable aleatoria definida como el tiempo que pasa, en horas, para que un radar detecte entre conductores sucesivos a los que exceden los límites de velocidad. La variable aleatoria X toma todos los valores de x para los que $x \geq 0$. ■

Definición 3.2: Si un espacio muestral contiene un número finito de posibilidades, o una serie interminable con tantos elementos como números enteros existen, se llama **espacio muestral discreto**.

Los resultados de algunos experimentos estadísticos no pueden ser ni finitos ni contables. Éste es el caso, por ejemplo, en una investigación que se realiza para medir las distancias que recorre un automóvil de cierta marca, en una ruta de prueba preestablecida, con cinco litros de gasolina. Si se asume que la distancia es una variable que se mide con algún grado de precisión, entonces salta a la vista que tenemos un número infinito de distancias posibles en el espacio muestral, que no se pueden igualar a la cantidad total de números enteros. Lo mismo ocurre en el caso de un experimento en que se registra el tiempo requerido para que ocurra una reacción química, en donde una vez más los posibles intervalos de tiempo que forman el espacio muestral serían un número infinito e incontable. Vemos ahora que no todos los espacios muestrales necesitan ser discretos.

Definición 3.3: Si un espacio muestral contiene un número infinito de posibilidades, igual al número de puntos en un segmento de recta, se le denomina **espacio muestral continuo**.

Una variable aleatoria se llama **variable aleatoria discreta** si se puede contar su conjunto de resultados posibles. En los ejemplos 3.1 a 3.5 las variables aleatorias son discretas. Sin embargo, una variable aleatoria cuyo conjunto de valores posibles es un intervalo completo de números no es discreta. Cuando una variable aleatoria puede tomar valores

en una escala continua, se le denomina **variable aleatoria continua**. A menudo los posibles valores de una variable aleatoria continua son precisamente los mismos valores incluidos en el espacio muestral continuo. Es evidente que las variables aleatorias descritas en los ejemplos 3.6 y 3.7 son variables aleatorias continuas.

En la mayoría de los problemas prácticos las variables aleatorias continuas representan datos *medidos*, como serían todos los posibles pesos, alturas, temperaturas, distancias o periodos de vida; en tanto que las variables aleatorias discretas representan datos *por conteo*, como el número de artículos defectuosos en una muestra de k artículos o el número de accidentes de carretera por año en una entidad específica. Observe que tanto Y como M , las variables aleatorias de los ejemplos 3.1 y 3.2, representan datos por conteo: Y el número de bolas rojas y M el número de emparejamientos correctos de cascos.

3.2 Distribuciones discretas de probabilidad

Una variable aleatoria discreta toma cada uno de sus valores con cierta probabilidad. Al lanzar una moneda tres veces, la variable X , que representa el número de caras, toma el valor 2 con $3/8$ de probabilidad, pues 3 de los 8 puntos muestrales igualmente probables tienen como resultado dos caras y una cruz. Si se suponen pesos iguales para los eventos simples del ejemplo 3.2, la probabilidad de que ningún obrero reciba el casco correcto, es decir, la probabilidad de que M tome el valor cero, es $1/3$. Los valores posibles m de M y sus probabilidades son

m	0	1	3
$P(M = m)$	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{6}$

Observe que los valores de m agotan todos los casos posibles, por lo tanto, las probabilidades suman 1.

Con frecuencia es conveniente representar todas las probabilidades de una variable aleatoria X usando una fórmula, la cual necesariamente sería una función de los valores numéricos x que denotaremos con $f(x)$, $g(x)$, $r(x)$ y así sucesivamente. Por lo tanto, escribimos $f(x) = P(X = x)$; es decir, $f(3) = P(X = 3)$. Al conjunto de pares ordenados $(x, f(x))$ se le llama **función de probabilidad**, **función de masa de probabilidad** o **distribución de probabilidad** de la variable aleatoria discreta X .

Definición 3.4: El conjunto de pares ordenados $(x, f(x))$ es una **función de probabilidad**, una **función de masa de probabilidad** o una **distribución de probabilidad** de la variable aleatoria discreta X si, para cada resultado x posible,

1. $f(x) \geq 0$,
2. $\sum_x f(x) = 1$,
3. $P(X = x) = f(x)$.

Ejemplo 3.8: Un embarque de 20 computadoras portátiles similares para una tienda minorista contiene 3 que están defectuosas. Si una escuela compra al azar 2 de estas computadoras, calcule la distribución de probabilidad para el número de computadoras defectuosas.

Solución: Sea X una variable aleatoria cuyos valores x son los números posibles de computadoras defectuosas compradas por la escuela. Entonces x sólo puede asumir los números 0, 1 y 2. Así,

$$f(0) = P(X = 0) = \frac{\binom{3}{0}\binom{17}{2}}{\binom{20}{2}} = \frac{68}{95}, \quad f(1) = P(X = 1) = \frac{\binom{3}{1}\binom{17}{1}}{\binom{20}{2}} = \frac{51}{190},$$

$$f(2) = P(X = 2) = \frac{\binom{3}{2}\binom{17}{0}}{\binom{20}{2}} = \frac{3}{190}.$$

Por consiguiente, la distribución de probabilidad de X es

x	0	1	2
$f(x)$	$\frac{68}{95}$	$\frac{51}{190}$	$\frac{3}{190}$

Ejemplo 3.9: Si una agencia automotriz vende 50% de su inventario de cierto vehículo extranjero equipado con bolsas de aire laterales, calcule una fórmula para la distribución de probabilidad del número de automóviles con bolsas de aire laterales entre los siguientes 4 vehículos que venda la agencia.

Solución: Como la probabilidad de vender un automóvil con bolsas de aire laterales es 0.5, los $2^4 = 16$ puntos del espacio muestral tienen la misma probabilidad de ocurrencia. Por lo tanto, el denominador para todas las probabilidades, y también para nuestra función, es 16. Para obtener el número de formas de vender tres automóviles con bolsas de aire laterales necesitamos considerar el número de formas de dividir 4 resultados en 2 celdas, con 3 automóviles con bolsas de aire laterales asignados a una celda, y el modelo sin bolsas de aire laterales asignado a la otra. Esto se puede hacer de $\binom{4}{3} = 4$ formas. En general, el evento de vender x modelos con bolsas de aire laterales y $4 - x$ modelos sin bolsas de aire laterales puede ocurrir de $\binom{4}{x}$ formas, donde x puede ser 0, 1, 2, 3 o 4. Por consiguiente, la distribución de probabilidad $f(x) = P(X = x)$ es

$$f(x) = \frac{1}{16} \binom{4}{x}, \quad \text{para } x = 0, 1, 2, 3, 4.$$

Existen muchos problemas en los que desearíamos calcular la probabilidad de que el valor observado de una variable aleatoria X sea menor o igual que algún número real x . Al escribir $F(x) = P(X \leq x)$ para cualquier número real x , definimos $F(x)$ como la **función de la distribución acumulativa** de la variable aleatoria X .

Definición 3.5: La **función de la distribución acumulativa** $F(x)$ de una variable aleatoria discreta X con distribución de probabilidad $f(x)$ es

$$F(x) = P(X \leq x) = \sum_{t \leq x} f(t), \quad \text{para } -\infty < x < \infty.$$

Para la variable aleatoria M , el número de emparejamientos correctos en el ejemplo 3.2, tenemos

$$F(2) = P(M \leq 2) = f(0) + f(1) = \frac{1}{3} + \frac{1}{2} = \frac{5}{6}.$$

La función de la distribución acumulativa de M es

$$F(m) = \begin{cases} 0, & \text{para } m < 0, \\ \frac{1}{3}, & \text{para } 0 \leq m < 1, \\ \frac{5}{6}, & \text{para } 1 \leq m < 3, \\ 1, & \text{para } m \geq 3. \end{cases}$$

Es necesario observar en particular el hecho de que la función de la distribución acumulativa es una función no decreciente monótona, la cual no sólo se define para los valores que toma la variable aleatoria dada sino para todos los números reales.

Ejemplo 3.10: Calcule la función de la distribución acumulativa de la variable aleatoria X del ejemplo 3.9. Utilice $F(x)$ para verificar que $f(2) = 3/8$.

Solución: El cálculo directo de la distribución de probabilidad del ejemplo 3.9 da $f(0) = 1/16$, $f(1) = 1/4$, $f(2) = 3/8$, $f(3) = 1/4$ y $f(4) = 1/16$. Por lo tanto,

$$F(0) = f(0) = \frac{1}{16},$$

$$F(1) = f(0) + f(1) = \frac{5}{16},$$

$$F(2) = f(0) + f(1) + f(2) = \frac{11}{16},$$

$$F(3) = f(0) + f(1) + f(2) + f(3) = \frac{15}{16},$$

$$F(4) = f(0) + f(1) + f(2) + f(3) + f(4) = 1.$$

Por lo tanto,

$$F(x) = \begin{cases} 0, & \text{para } x < 0, \\ \frac{1}{16}, & \text{para } 0 \leq x < 1, \\ \frac{5}{16}, & \text{para } 1 \leq x < 2, \\ \frac{11}{16}, & \text{para } 2 \leq x < 3, \\ \frac{15}{16}, & \text{para } 3 \leq x < 4, \\ 1 & \text{para } x \geq 4. \end{cases}$$

Entonces,

$$f(2) = F(2) - F(1) = \frac{11}{16} - \frac{5}{16} = \frac{3}{8}.$$

A menudo es útil ver una distribución de probabilidad en forma gráfica. Se pueden graficar los puntos $(x, f(x))$ del ejemplo 3.9 para obtener la figura 3.1. Si unimos los puntos al eje x , ya sea con una línea punteada o con una línea sólida, obtenemos una gráfica de función de masa de probabilidad. La figura 3.1 permite ver fácilmente qué valores de X tienen más probabilidad de ocurrencia y, en este caso, también indica una situación perfectamente simétrica.

Sin embargo, en vez de graficar los puntos $(x, f(x))$, lo que hacemos más a menudo es construir rectángulos como en la figura 3.2. Aquí los rectángulos se construyen de manera que sus bases, con la misma anchura, se centren en cada valor x , y que sus alturas iguallen a las probabilidades correspondientes dadas por $f(x)$. Las bases se construyen de forma tal que no dejen espacios entre los rectángulos. La figura 3.2 se denomina **histograma de probabilidad**.

Como cada base en la figura 3.2 tiene el ancho de una unidad, $P(X = x)$ es igual al área del rectángulo centrado en x . Incluso si las bases no tuvieran el ancho de una unidad, podríamos ajustar las alturas de los rectángulos para que tengan áreas que iguallen las probabilidades de X de tomar cualquiera de sus valores x . Este concepto de utilizar

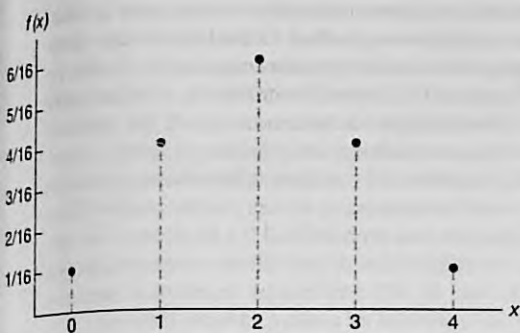


Figura 3.1: Gráfica de función de masa de probabilidad.

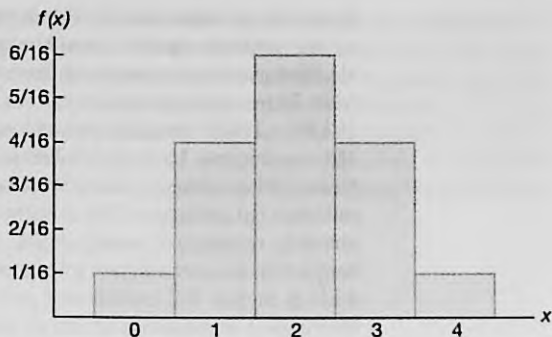


Figura 3.2: Histograma de probabilidad.

áreas para representar probabilidades es necesario para nuestro estudio de la distribución de probabilidad de una variable aleatoria continua.

La gráfica de la función de la distribución acumulativa del ejemplo 3.9, que aparece como una función escalonada en la figura 3.3, se obtiene graficando los puntos $(x, F(x))$.

Ciertas distribuciones de probabilidad se aplican a más de una situación física. La distribución de probabilidad del ejemplo 3.9 también se aplica a la variable aleatoria Y , donde Y es el número de caras que se obtienen cuando una moneda se lanza 4 veces, o a la variable aleatoria W , donde W es el número de cartas rojas que resultan cuando se sacan 4 cartas al azar de una baraja de manera sucesiva, se reemplaza cada carta y se baraja antes de sacar la siguiente. En el capítulo 5 se estudiarán distribuciones discretas especiales que se aplican a diversas situaciones experimentales.

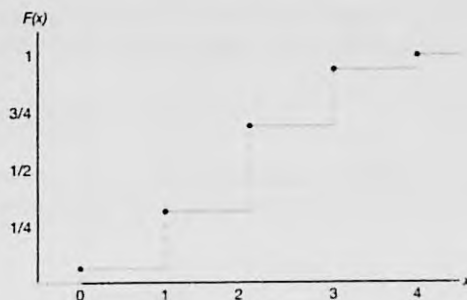


Figura 3.3: Función de distribución acumulativa discreta.

3.3 Distribuciones de probabilidad continua

Una variable aleatoria continua tiene una probabilidad 0 de adoptar *exactamente* cualquiera de sus valores. En consecuencia, su distribución de probabilidad no se puede

presentar en forma tabular. En un principio esto parecería sorprendente, pero se vuelve más probable cuando consideramos un ejemplo específico. Consideremos una variable aleatoria cuyos valores son las estaturas de todas las personas mayores de 21 años de edad. Entre cualesquiera dos valores, digamos 163.5 y 164.5 centímetros, o incluso entre 163.99 y 164.01 centímetros, hay un número infinito de estaturas, una de las cuales es 164 centímetros. La probabilidad de seleccionar al azar a una persona que tenga exactamente 164 centímetros de estatura en lugar de una del conjunto infinitamente grande de estaturas tan cercanas a 164 centímetros que humanamente no sea posible medir la diferencia es remota, por consiguiente, asignamos una probabilidad 0 a tal evento. Sin embargo, esto no ocurre si nos referimos a la probabilidad de seleccionar a una persona que mida al menos 163 centímetros pero no más de 165 centímetros de estatura. Aquí nos referimos a un intervalo en vez de a un valor puntual de nuestra variable aleatoria.

Nos interesamos por el cálculo de probabilidades para varios intervalos de variables aleatorias continuas como $P(a < X < b)$, $P(W \geq c)$, etc. Observe que cuando X es continua,

$$P(a < X \leq b) = P(a < X < b) + P(X = b) = P(a < X < b).$$

Es decir, no importa si incluimos o no un extremo del intervalo. Sin embargo, esto no es cierto cuando X es discreta.

Aunque la distribución de probabilidad de una variable aleatoria continua no se puede representar de forma tabular, sí es posible plantearla como una fórmula, la cual necesariamente será función de los valores numéricos de la variable aleatoria continua X , y como tal se representará mediante la notación funcional $f(x)$. Cuando se trata con variables continuas, a $f(x)$ por lo general se le llama **función de densidad de probabilidad**, o simplemente **función de densidad** de X . Como X se define sobre un espacio muestral continuo, es posible que $f(x)$ tenga un número finito de discontinuidades. Sin embargo, la mayoría de las funciones de densidad que tienen aplicaciones prácticas en el análisis de datos estadísticos son continuas y sus gráficas pueden tomar cualesquiera de varias formas, algunas de las cuales se presentan en la figura 3.4. Como se utilizarán áreas para representar probabilidades y éstas son valores numéricos positivos, la función de densidad debe caer completamente arriba del eje x .

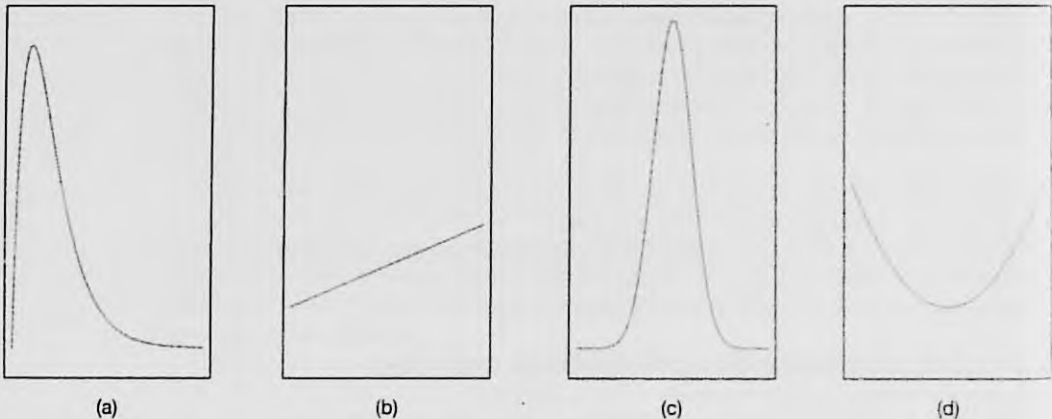


Figura 3.4: Funciones de densidad típicas.

Una función de densidad de probabilidad se construye de manera que el área bajo su curva limitada por el eje x sea igual a 1, cuando se calcula en el rango de X para el que se define $f(x)$. Como este rango de X es un intervalo finito, siempre es posible extender el intervalo para que incluya a todo el conjunto de números reales definiendo $f(x)$ como cero en todos los puntos de las partes extendidas del intervalo. En la figura 3.5 la probabilidad de que X tome un valor entre a y b es igual al área sombreada bajo la función de densidad entre las ordenadas en $x = a$ y $x = b$, y a partir del cálculo integral está dada por

$$P(a < X < b) = \int_a^b f(x) dx.$$

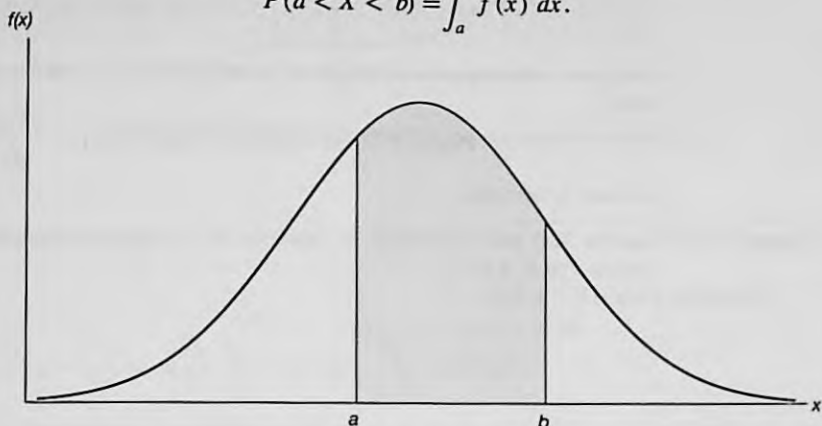


Figura 3.5: $P(a < X < b)$.

Definición 3.6: La función $f(x)$ es una **función de densidad de probabilidad** (fdp) para la variable aleatoria continua X , definida en el conjunto de números reales, si

1. $f(x) \geq 0$, para toda $x \in R$.
2. $\int_{-\infty}^{\infty} f(x) dx = 1$.
3. $P(a < X < b) = \int_a^b f(x) dx$.

Ejemplo 3.11: Suponga que el error en la temperatura de reacción, en $^{\circ}\text{C}$, en un experimento de laboratorio controlado, es una variable aleatoria continua X que tiene la función de densidad de probabilidad

$$f(x) = \begin{cases} \frac{x^2}{3}, & -1 < x < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Verifique que $f(x)$ es una función de densidad.
- b) Calcule $P(0 < X \leq 1)$.

Solución: Usamos la definición 3.6.

- a) Evidentemente, $f(x) \geq 0$. Para verificar la condición 2 de la definición 3.6 tenemos

$$\int_{-\infty}^{\infty} f(x) dx = \int_{-1}^2 \frac{x^2}{3} dx = \frac{x^3}{9} \Big|_{-1}^2 = \frac{8}{9} + \frac{1}{9} = 1.$$

b) Si usamos la fórmula 3 de la definición 3.6, obtenemos

$$P(0 < X \leq 1) = \int_0^1 \frac{x^2}{3} dx = \frac{x^3}{9} \Big|_0^1 = \frac{1}{9}.$$

Definición 3.7: La función de distribución acumulativa $F(x)$, de una variable aleatoria continua X con función de densidad $f(x)$, es

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt, \text{ para } -\infty < x < \infty.$$

Como una consecuencia inmediata de la definición 3.7 se pueden escribir los dos resultados,

$$P(a < X < b) = F(b) - F(a) \text{ y } f(x) = \frac{dF(x)}{dx},$$

si existe la derivada.

Ejemplo 3.12: Calcule $F(x)$ para la función de densidad del ejemplo 3.11 y utilice el resultado para evaluar $P(0 < X \leq 1)$.

Solución: Para $-1 < x < 2$,

$$F(x) = \int_{-\infty}^x f(t) dt = \int_{-1}^x \frac{t^2}{3} dt = \frac{t^3}{9} \Big|_{-1}^x = \frac{x^3 + 1}{9}.$$

Por lo tanto,

$$F(x) = \begin{cases} 0, & x < -1, \\ \frac{x^3+1}{9}, & -1 \leq x < 2, \\ 1, & x \geq 2. \end{cases}$$

La función de la distribución acumulativa $F(x)$ se expresa en la figura 3.6. Entonces,

$$P(0 < X \leq 1) = F(1) - F(0) = \frac{2}{9} - \frac{1}{9} = \frac{1}{9},$$

que coincide con el resultado que se obtuvo al utilizar la función de densidad en el ejemplo 3.11.

Ejemplo 3.13: El Departamento de Energía (DE) asigna proyectos mediante licitación y, por lo general, estima lo que debería ser una licitación razonable. Sea b el estimado. El DE determinó que la función de densidad de la licitación ganadora (baja) es

$$f(y) = \begin{cases} \frac{5}{8b}, & \frac{2}{5}b \leq y \leq 2b, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule $F(y)$ y utilice el resultado para determinar la probabilidad de que la licitación ganadora sea menor que la estimación preliminar b del DE.

Solución: Para $2b/5 \leq y \leq 2b$,

$$F(y) = \int_{2b/5}^y \frac{5}{8b} dy = \frac{5y}{8b} \Big|_{2b/5}^y = \frac{5y}{8b} - \frac{1}{4}.$$

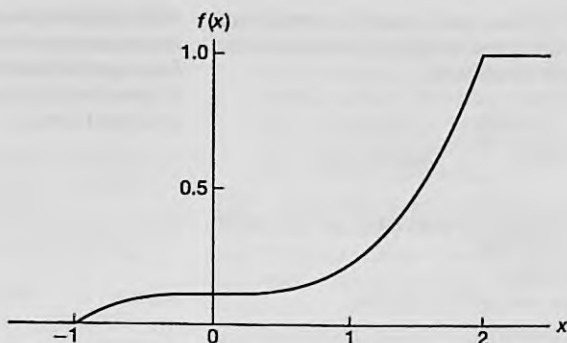


Figura 3.6: Función de distribución acumulativa continua.

Por consiguiente,

$$F(y) = \begin{cases} 0, & y < \frac{2}{5}b, \\ \frac{5y}{8b} - \frac{1}{4}, & \frac{2}{5}b \leq y < 2b, \\ 1, & y \geq 2b. \end{cases}$$

Para determinar la probabilidad de que la licitación ganadora sea menor que la estimación preliminar b de la licitación tenemos

$$P(Y \leq b) = F(b) = \frac{5}{8} - \frac{1}{4} = \frac{3}{8}.$$

■

Ejercicios

3.1 Clasifique las siguientes variables aleatorias como discretas o continuas:

X : el número de accidentes automovilísticos que ocurren al año en Virginia.

Y : el tiempo para jugar 18 hoyos de golf.

M : la cantidad de leche que una vaca específica produce anualmente.

N : el número de huevos que una gallina pone mensualmente.

P : el número de permisos para construcción que los funcionarios de una ciudad emiten cada mes.

Q : el peso del grano producido por acre.

3.2 Un embarque foráneo de 5 automóviles extranjeros contiene 2 que tienen ligeras manchas de pintura. Suponga que una agencia recibe 3 de estos automóviles al azar y liste los elementos del espacio muestral S usando las letras M y N para “manchado” y “sin mancha”, respectivamente; luego asigne a cada punto

muestral un valor x de la variable aleatoria X que representa el número de automóviles con manchas de pintura que compró la agencia.

3.3 Sea W la variable aleatoria que da el número de caras menos el número de cruces en tres lanzamientos de una moneda. Liste los elementos del espacio muestral S para los tres lanzamientos de la moneda y asigne un valor w de W a cada punto muestral.

3.4 Se lanza una moneda hasta que se presentan 3 caras sucesivamente. Liste sólo aquellos elementos del espacio muestral que requieren 6 o menos lanzamientos. ¿Es éste un espacio muestral discreto? Explique su respuesta.

3.5 Determine el valor c de modo que cada una de las siguientes funciones sirva como distribución de probabilidad de la variable aleatoria discreta X :

a) $f(x) = c(x^2 + 4)$, para $x = 0, 1, 2, 3$;

b) $f(x) = c \binom{2}{x} \binom{3}{3-x}$, para $x = 0, 1, 2$.

3.6 La vida útil, en días, para frascos de cierta medicina de prescripción es una variable aleatoria que tiene la siguiente función de densidad:

$$f(x) = \begin{cases} \frac{20,000}{(x+100)^3}, & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la probabilidad de que un frasco de esta medicina tenga una vida útil de

- al menos 200 días;
- cualquier lapso entre 80 y 120 días.

3.7 El número total de horas, medidas en unidades de 100 horas, que una familia utiliza una aspiradora en un periodo de un año es una variable aleatoria continua X que tiene la siguiente función de densidad:

$$f(x) = \begin{cases} x, & 0 < x < 1, \\ 2 - x, & 1 \leq x < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la probabilidad de que en un periodo de un año una familia utilice su aspiradora

- menos de 120 horas;
- entre 50 y 100 horas.

3.8 Obtenga la distribución de probabilidad de la variable aleatoria W del ejercicio 3.3; suponga que la moneda está cargada, de manera que existe el doble de probabilidad de que ocurra una cara que una cruz.

3.9 La proporción de personas que responden a cierta encuesta enviada por correo es una variable aleatoria continua X que tiene la siguiente función de densidad:

$$f(x) = \begin{cases} \frac{2(x+2)}{5}, & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Demuestre que $P(0 < X < 1) = 1$.
- Calcule la probabilidad de que más de $1/4$ pero menos de $1/2$ de las personas contactadas respondan a este tipo de encuesta.

3.10 Encuentre una fórmula para la distribución de probabilidad de la variable aleatoria X que represente el resultado cuando se lanza un dado una vez.

3.11 Un embarque de 7 televisores contiene 2 unidades defectuosas. Un hotel compra 3 de los televisores al azar. Si x es el número de unidades defectuosas que compra el hotel, calcule la distribución de probabilidad de X . Expresé los resultados de forma gráfica como un histograma de probabilidad.

3.12 Una empresa de inversiones ofrece a sus clientes bonos municipales que vencen después de varios años. Dado que la función de distribución acumulativa de T , el número de años para el vencimiento de un bono que se elige al azar, es

$$F(t) = \begin{cases} 0, & t < 1, \\ \frac{1}{4}, & 1 \leq t < 3, \\ \frac{1}{2}, & 3 \leq t < 5, \\ \frac{3}{4}, & 5 \leq t < 7, \\ 1, & t \geq 7, \end{cases}$$

calcule

- $P(T = 5)$;
- $P(T > 3)$;
- $P(1.4 < T < 6)$;
- $P(T \leq 5 \mid T \geq 2)$;

3.13 La distribución de probabilidad de X , el número de imperfecciones que se encuentran en cada 10 metros de una tela sintética que viene en rollos continuos de ancho uniforme, está dada por

x	0	1	2	3	4
$f(x)$	0.41	0.37	0.16	0.05	0.01

Construya la función de distribución acumulativa de X .

3.14 El tiempo que pasa, en horas, para que un radar detecte entre conductores sucesivos a los que exceden los límites de velocidad es una variable aleatoria continua con una función de distribución acumulativa

$$F(x) = \begin{cases} 0, & x < 0, \\ 1 - e^{-8x}, & x \geq 0. \end{cases}$$

Calcule la probabilidad de que el tiempo que pase para que el radar detecte entre conductores sucesivos a los que exceden los límites de velocidad sea menor de 12 minutos

- usando la función de distribución acumulativa de X ;
- utilizando la función de densidad de probabilidad de X .

3.15 Calcule la función de distribución acumulativa de la variable aleatoria X que represente el número de unidades defectuosas en el ejercicio 3.11. Luego, utilice $F(x)$ para calcular

- $P(X = 1)$;
- $P(0 < X \leq 2)$.

3.16 Construya una gráfica de la función de distribución acumulativa del ejercicio 3.15.

3.17 Una variable aleatoria continua X , que puede tomar valores entre $x = 1$ y $x = 3$, tiene una función de densidad dada por $f(x) = 1/2$.

- Muestre que el área bajo la curva es igual a 1.
- Calcule $P(2 < X < 2.5)$.
- Calcule $P(X \leq 1.6)$.

3.18 Una variable aleatoria continua X , que puede tomar valores entre $x = 2$ y $x = 5$, tiene una función de densidad dada por $f(x) = 2(1 + x)/27$. Calcule

- $P(X < 4)$;
- $P(3 \leq X < 4)$.

3.19 Para la función de densidad del ejercicio 3.17 calcule $F(x)$. Utilícela para evaluar $P(2 < X < 2.5)$.

3.20 Para la función de densidad del ejercicio 3.18 calcule $F(x)$ y utilícela para evaluar $P(3 \leq X < 4)$.

3.21 Considere la función de densidad

$$f(x) = \begin{cases} k\sqrt{x}, & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Evalúe k .
- Calcule $F(x)$ y utilice el resultado para evaluar

$$P(0.3 < X < 0.6).$$

3.22 Se sacan tres cartas de una baraja de manera sucesiva y sin reemplazo. Calcule la distribución de probabilidad para la cantidad de espadas.

3.23 Calcule la función de distribución acumulativa de la variable aleatoria W del ejercicio 3.8. Use $F(w)$ para calcular

- $P(W > 0)$;
- $P(-1 \leq W < 3)$.

3.24 Calcule la distribución de probabilidad para el número de discos compactos de jazz cuando, de una colección que consta de 5 de jazz, 2 de música clásica y 3 de rock, se seleccionan 4 CD al azar. Expresé sus resultados utilizando una fórmula.

3.25 De una caja que contiene 4 monedas de 10 centavos y 2 monedas de 5 centavos se seleccionan 3 monedas al azar y sin reemplazo. Calcule la distribución de probabilidad para el total T de las 3 monedas. Expresé la distribución de probabilidad de forma gráfica como un histograma de probabilidad.

3.26 De una caja que contiene 4 bolas negras y 2 verdes se sacan 3 bolas sucesivamente, cada bola se regresa a la caja antes de sacar la siguiente. Calcule la distribución de probabilidad para el número de bolas verdes.

3.27 El tiempo que pasa, en horas, antes de que una parte importante de un equipo electrónico que se utiliza para fabricar un reproductor de DVD empiece a fallar tiene la siguiente función de densidad:

$$f(x) = \begin{cases} \frac{1}{2000} \exp(-x/2000), & x \geq 0, \\ 0, & x < 0. \end{cases}$$

- Calcule $F(x)$.
- Determine la probabilidad de que el componente (y, por lo tanto, el reproductor de DVD) funcione durante más de 1000 horas antes de que sea necesario reemplazar el componente.
- Determine la probabilidad de que el componente falle antes de 2000 horas.

3.28 Un productor de cereales está consciente de que el peso del producto varía ligeramente entre una y otra caja. De hecho, cuenta con suficientes datos históricos para determinar la función de densidad que describe la estructura de probabilidad para el peso (en onzas). Si X es el peso, en onzas, de la variable aleatoria, la función de densidad se describe como

$$f(x) = \begin{cases} \frac{2}{5}, & 23.75 \leq x \leq 26.25, \\ 0, & \text{en otro caso.} \end{cases}$$

- Verifique que sea una función de densidad válida.
- Determine la probabilidad de que el peso sea menor que 24 onzas.
- La empresa desea que un peso mayor que 26 onzas sea un caso extraordinariamente raro. ¿Cuál será la probabilidad de que en verdad ocurra este caso extraordinariamente raro?

3.29 Un factor importante en el combustible sólido para proyectiles es la distribución del tamaño de las partículas. Cuando las partículas son demasiado grandes se presentan problemas importantes. A partir de datos de producción históricos se determinó que la distribución del tamaño (en micras) de las partículas se caracteriza por

$$f(x) = \begin{cases} 3x^{-4}, & x > 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Verifique que sea una función de densidad válida.
- Evalúe $F(x)$.
- ¿Cuál es la probabilidad de que una partícula tomada al azar del combustible fabricado sea mayor que 4 micras?

3.30 Las mediciones en los sistemas científicos siempre están sujetas a variación, algunas veces más que otras. Hay muchas estructuras para los errores de medición y los estadísticos pasan mucho tiempo modelándolos. Suponga que el error de medición X de cierta cantidad física es determinado por la siguiente función de densidad:

$$f(x) = \begin{cases} k(3 - x^2), & -1 \leq x \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Determine k , que representa $f(x)$, una función de densidad válida.
- Calcule la probabilidad de que un error aleatorio en la medición sea menor que $\frac{1}{2}$.
- Para esta medición específica, resulta indeseable si la magnitud del error (es decir, $|x|$) es mayor que 0.8. ¿Cuál es la probabilidad de que esto ocurra?

3.31 Con base en pruebas extensas, el fabricante de una lavadora determinó que el tiempo Y (en años) para que el electrodoméstico requiera una reparación mayor se obtiene mediante la siguiente función de densidad de probabilidad:

$$f(y) = \begin{cases} \frac{1}{4}e^{-y/4}, & y \geq 0, \\ 0, & \text{en cualquier otro caso.} \end{cases}$$

- Los críticos considerarían que la lavadora es una ganga si no hay probabilidades de que requiera una reparación mayor antes del sexto año. Comente sobre esto determinando $P(Y > 6)$.
- ¿Cuál es la probabilidad de que la lavadora requiera una reparación mayor durante el primer año?

3.32 Se está revisando qué proporciones de su presupuesto asigna cierta empresa industrial a controles ambientales y de contaminación. Un proyecto de recopilación de datos determina que la distribución de tales proporciones está dada por

$$f(y) = \begin{cases} 5(1-y)^4, & 0 \leq y \leq 1, \\ 0, & \text{en cualquier otro caso.} \end{cases}$$

- Verifique que la función de densidad anterior sea válida.
- ¿Cuál es la probabilidad de que una empresa elegida al azar gaste menos de 10% de su presupuesto en controles ambientales y de contaminación?
- ¿Cuál es la probabilidad de que una empresa seleccionada al azar gaste más de 50% de su presupuesto en controles ambientales y de la contaminación?

3.33 Suponga que cierto tipo de pequeñas empresas de procesamiento de datos están tan especializadas que algunas tienen dificultades para obtener utilidades durante su primer año de operación. La función de densidad de probabilidad que caracteriza la proporción Y que obtiene utilidades está dada por

$$f(y) = \begin{cases} ky^4(1-y)^3, & 0 \leq y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- ¿Cuál es el valor de k que hace de la anterior una función de densidad válida?
- Calcule la probabilidad de que al menos 50% de las empresas tenga utilidades durante el primer año.
- Calcule la probabilidad de que al menos 80% de las empresas tenga utilidades durante el primer año.

3.34 Los tubos de magnetrón se producen en una línea de ensamble automatizada. Periódicamente se utiliza un plan de muestreo para evaluar la calidad en la longitud de los tubos; sin embargo, dicha medida está sujeta a incertidumbre. Se considera que la probabilidad de que un tubo elegido al azar cumpla con las especificaciones de longitud es 0.99. Se utiliza un plan de muestreo en el cual se mide la longitud de 5 tubos elegidos al azar.

- Muestre que la función de probabilidad de Y , el número de tubos de cada 5 que cumplen con las especificaciones de longitud, está dada por la siguiente función de probabilidad discreta:

$$f(y) = \frac{5!}{y!(5-y)!} (0.99)^y (0.01)^{5-y},$$

- Suponga que se eligen artículos de la línea al azar y 3 no cumplen con las especificaciones. Utilice la $f(y)$ anterior para apoyar o refutar la conjetura de que hay 0.99 de probabilidades de que un solo tubo cumpla con las especificaciones.

3.35 Suponga que a partir de gran cantidad de datos históricos se sabe que X , el número de automóviles que llegan a una intersección específica durante un periodo de 20 segundos, se determina mediante la siguiente función de probabilidad discreta

$$f(x) = e^{-6} \frac{6^x}{x!}, \quad \text{para } x = 0, 1, 2, \dots$$

- Calcule la probabilidad de que en un periodo específico de 20 segundos más de 8 automóviles lleguen a la intersección.
- Calcule la probabilidad de que sólo lleguen 2 automóviles.

3.36 En una tarea de laboratorio, si el equipo está funcionando, la función de densidad del resultado observado, X , es

$$f(x) = \begin{cases} 2(1-x), & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Calcule $P(X \leq 1/3)$.
- ¿Cuál es la probabilidad de que X sea mayor que 0.5?
- Dado que $X \geq 0.5$, ¿cuál es la probabilidad de que X sea menor que 0.75?

3.4 Distribuciones de probabilidad conjunta

El estudio de las variables aleatorias y sus distribuciones de probabilidad de la sección anterior se restringió a espacios muestrales unidimensionales, ya que registramos los resultados de un experimento como los valores que toma una sola variable aleatoria. No

obstante, habrá situaciones en las que se busque registrar los resultados simultáneos de diversas variables aleatorias. Por ejemplo, en un experimento químico controlado podríamos medir la cantidad del precipitado P y la del volumen V de gas liberado, lo que daría lugar a un espacio muestral bidimensional que consta de los resultados (p, v) ; o bien, podríamos interesarnos en la dureza d y en la resistencia a la tensión T de cobre estirado en frío que produciría los resultados (d, t) . En un estudio realizado con estudiantes universitarios para determinar la probabilidad de que tengan éxito en la universidad, basado en los datos del nivel preparatoria, se podría utilizar un espacio muestral tridimensional y registrar la calificación que obtuvo cada uno en la prueba de aptitudes, el lugar que cada uno ocupó en la preparatoria y la calificación promedio que cada uno obtuvo al final de su primer año en la universidad.

Si X y Y son dos variables aleatorias discretas, la distribución de probabilidad para sus ocurrencias simultáneas se representa mediante una función con valores $f(x, y)$, para cualquier par de valores (x, y) dentro del rango de las variables aleatorias X y Y . Se acostumbra referirse a esta función como la **distribución de probabilidad conjunta** de X y Y .

Por consiguiente, en el caso discreto,

$$f(x, y) = P(X = x, Y = y);$$

es decir, los valores $f(x, y)$ dan la probabilidad de que los resultados x y y ocurran al mismo tiempo. Por ejemplo, si se le va a dar servicio a los neumáticos de un camión de transporte pesado, y X representa el número de millas que éstos han recorrido y Y el número de neumáticos que deben ser reemplazados, entonces $f(30,000, 5)$ es la probabilidad de que los neumáticos hayan recorrido más de 30,000 millas y que el camión necesite 5 neumáticos nuevos.

Definición 3.8: La función $f(x, y)$ es una **distribución de probabilidad conjunta** o **función de masa de probabilidad** de las variables aleatorias discretas X y Y , si

1. $f(x, y) \geq 0$ para toda (x, y) ,
2. $\sum_x \sum_y f(x, y) = 1$,
3. $P(X = x, Y = y) = f(x, y)$.

Para cualquier región A en el plano xy , $P[(X, Y) \in A] = \sum_A f(x, y)$.

Ejemplo 3.14: Se seleccionan al azar 2 repuestos para bolígrafo de una caja que contiene 3 repuestos azules, 2 rojos y 3 verdes. Si X es el número de repuestos azules y Y es el número de repuestos rojos seleccionados, calcule

- a) la función de probabilidad conjunta $f(x, y)$,
- b) $P[(X, Y) \in A]$, donde A es la región $\{(x, y) | x + y \leq 1\}$.

Solución: Los posibles pares de valores (x, y) son $(0, 0)$, $(0, 1)$, $(1, 0)$, $(1, 1)$, $(0, 2)$ y $(2, 0)$.

- a) Ahora bien, $f(0, 1)$, por ejemplo, representa la probabilidad de seleccionar un repuesto rojo y uno verde. El número total de formas igualmente probables de seleccionar cualesquiera 2 repuestos de los 8 es $\binom{8}{2} = 28$. El número de formas de seleccionar 1 rojo de 2 repuestos rojos y 1 verde de 3 repuestos verdes es $\binom{2}{1} \binom{3}{1} = 6$. En consecuencia, $f(0, 1) = 6/28 = 3/14$. Cálculos similares dan las probabilidades para

los otros casos, los cuales se presentan en la tabla 3.1. Observe que las probabilidades suman 1. En el capítulo 5 se volverá evidente que la distribución de probabilidad conjunta de la tabla 3.1 se puede representar con la fórmula

$$f(x, y) = \frac{\binom{3}{x} \binom{2}{y} \binom{2-x-y}{2}}{\binom{8}{2}},$$

para $x = 0, 1, 2; y = 0, 1, 2; y 0 \leq x + y \leq 2$.

b) La probabilidad de que (X, Y) caiga en la región A es

$$\begin{aligned} P[(X, Y) \in A] &= P(X + Y \leq 1) = f(0, 0) + f(0, 1) + f(1, 0) \\ &= \frac{3}{28} + \frac{3}{14} + \frac{9}{28} = \frac{9}{14}. \end{aligned}$$

Tabla 3.1: Distribución de probabilidad conjunta para el ejemplo 3.14

$f(x, y)$		x			Totales por renglón
		0	1	2	
y	0	$\frac{3}{28}$	$\frac{9}{28}$	$\frac{3}{28}$	$\frac{15}{28}$
	1	$\frac{3}{14}$	$\frac{3}{14}$	0	$\frac{3}{7}$
	2	$\frac{1}{28}$	0	0	$\frac{1}{28}$
Totales por columna		$\frac{5}{14}$	$\frac{15}{28}$	$\frac{3}{28}$	1

Cuando X y Y son variables aleatorias continuas, la **función de densidad conjunta** $f(x, y)$ es una superficie que yace sobre el plano xy , y $P[(X, Y) \in A]$, donde A es cualquier región en el plano xy , que es igual al volumen del cilindro recto limitado por la base A y la superficie.

Definición 3.9: La función $f(x, y)$ es una **función de densidad conjunta** de las variables aleatorias continuas X y Y si

1. $f(x, y) \geq 0$, para toda (x, y) ,
2. $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$,
3. $P[(X, Y) \in A] = \int \int_A f(x, y) dx dy$, para cualquier región A en el plano xy .

Ejemplo 3.15: Una empresa privada opera un local que da servicio a clientes que llegan en automóvil y otro que da servicio a clientes que llegan caminando. En un día elegido al azar, sean X y Y , respectivamente, las proporciones de tiempo que ambos locales están en servicio, y suponiendo que la función de densidad conjunta de estas variables aleatorias es

$$f(x, y) = \begin{cases} \frac{2}{5}(2x + 3y), & 0 \leq x \leq 1, 0 \leq y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

a) Verifique la condición 2 de la definición 3.9.

b) Calcule $P[(X, Y) \in A]$, donde $A = \{(x, y) \mid 0 < x < \frac{1}{2}, \frac{1}{4} < y < \frac{1}{2}\}$.

Solución: a) La integración de $f(x,y)$ sobre la totalidad de la región es

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x,y) dx dy &= \int_0^1 \int_0^1 \frac{2}{5}(2x+3y) dx dy \\ &= \int_0^1 \left(\frac{2x^2}{5} + \frac{6xy}{5} \right) \Big|_{x=0}^{x=1} dy \\ &= \int_0^1 \left(\frac{2}{5} + \frac{6y}{5} \right) dy = \left(\frac{2y}{5} + \frac{3y^2}{5} \right) \Big|_0^1 = \frac{2}{5} + \frac{3}{5} = 1. \end{aligned}$$

b) Para calcular la probabilidad utilizamos

$$\begin{aligned} P[(X, Y) \in A] &= P\left(0 < X < \frac{1}{2}, \frac{1}{4} < Y < \frac{1}{2}\right) \\ &= \int_{1/4}^{1/2} \int_0^{1/2} \frac{2}{5}(2x+3y) dx dy \\ &= \int_{1/4}^{1/2} \left(\frac{2x^2}{5} + \frac{6xy}{5} \right) \Big|_{x=0}^{x=1/2} dy = \int_{1/4}^{1/2} \left(\frac{1}{10} + \frac{3y}{5} \right) dy \\ &= \left(\frac{y}{10} + \frac{3y^2}{10} \right) \Big|_{1/4}^{1/2} \\ &= \frac{1}{10} \left[\left(\frac{1}{2} + \frac{3}{4} \right) - \left(\frac{1}{4} + \frac{3}{16} \right) \right] = \frac{13}{160}. \end{aligned}$$

Dada la distribución de probabilidad conjunta $f(x,y)$ de las variables aleatorias discretas X y Y , la distribución de probabilidad $g(x)$ sólo de X se obtiene sumando $f(x,y)$ sobre los valores de Y . De manera similar, la distribución de probabilidad $h(y)$ de sólo Y se obtiene sumando $f(x,y)$ sobre los valores de X . Definimos $g(x)$ y $h(y)$ como **distribuciones marginales** de X y Y , respectivamente. Cuando X y Y son variables aleatorias continuas, las sumatorias se reemplazan por integrales. Ahora podemos establecer la siguiente definición general.

Definición 3.10: Las **distribuciones marginales** sólo de X y sólo de Y son

$$g(x) = \sum_y f(x,y) \quad \text{y} \quad h(y) = \sum_x f(x,y)$$

para el caso discreto, y

$$g(x) = \int_{-\infty}^{\infty} f(x,y) dy \quad \text{y} \quad h(y) = \int_{-\infty}^{\infty} f(x,y) dx$$

para el caso continuo.

El término *marginal* se utiliza aquí porque, en el caso discreto, los valores de $g(x)$ y $h(y)$ son precisamente los totales marginales de las columnas y los renglones respectivos, cuando los valores de $f(x,y)$ se muestran en una tabla rectangular.

Ejemplo 3.16: Muestre que los totales de columnas y renglones de la tabla 3.1 dan las distribuciones marginales de sólo X y sólo Y .

Solución: Para la variable aleatoria X vemos que

$$g(0) = f(0, 0) + f(0, 1) + f(0, 2) = \frac{3}{28} + \frac{3}{14} + \frac{1}{28} = \frac{5}{14},$$

$$g(1) = f(1, 0) + f(1, 1) + f(1, 2) = \frac{9}{28} + \frac{3}{14} + 0 = \frac{15}{28},$$

y

$$g(2) = f(2, 0) + f(2, 1) + f(2, 2) = \frac{3}{28} + 0 + 0 = \frac{3}{28},$$

que son precisamente los totales por columna de la tabla 3.1. De manera similar podemos mostrar que los valores de $h(y)$ están dados por los totales de los renglones. En forma tabular, estas distribuciones marginales se pueden escribir como sigue:

x	0	1	2		y	0	1	2
$g(x)$	$\frac{5}{14}$	$\frac{15}{28}$	$\frac{3}{28}$		$h(y)$	$\frac{15}{28}$	$\frac{3}{7}$	$\frac{1}{28}$

Ejemplo 3.17: Calcule $g(x)$ y $h(y)$ para la función de densidad conjunta del ejemplo 3.15.

Solución: Por definición,

$$g(x) = \int_{-\infty}^{\infty} f(x, y) dy = \int_0^1 \frac{2}{5}(2x + 3y) dy = \left(\frac{4xy}{5} + \frac{6y^2}{10} \right) \Big|_{y=0}^{y=1} = \frac{4x + 3}{5},$$

para $0 \leq x \leq 1$, y $g(x) = 0$ en otro caso. De manera similar,

$$h(y) = \int_{-\infty}^{\infty} f(x, y) dx = \int_0^1 \frac{2}{5}(2x + 3y) dx = \frac{2(1 + 3y)}{5},$$

para $0 \leq y \leq 1$, y $h(y) = 0$ en otro caso.

El hecho de que las distribuciones marginales $g(x)$ y $h(y)$ sean en realidad las distribuciones de probabilidad de las variables individuales X y Y solas se puede verificar mostrando que se satisfacen las condiciones de la definición 3.4 o de la definición 3.6. Por ejemplo, en el caso continuo

$$\int_{-\infty}^{\infty} g(x) dx = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dy dx = 1,$$

y

$$P(a < X < b) = P(a < X < b, -\infty < Y < \infty)$$

$$= \int_a^b \int_{-\infty}^{\infty} f(x, y) dy dx = \int_a^b g(x) dx.$$

En la sección 3.1 establecimos que el valor x de la variable aleatoria X representa un evento que es un subconjunto del espacio muestral. Si utilizamos la definición de probabilidad condicional que se estableció en el capítulo 2,

$$P(B|A) = \frac{P(A \cap B)}{P(A)}, \text{ siempre que } P(A) > 0,$$

donde A y B son ahora los eventos definidos por $X = x$ y $Y = y$, respectivamente, entonces,

$$P(Y = y | X = x) = \frac{P(X = x, Y = y)}{P(X = x)} = \frac{f(x, y)}{g(x)}, \text{ siempre que } g(x) > 0,$$

donde X y Y son variables aleatorias discretas.

No es difícil mostrar que la función $f(x, y)/g(x)$, que es estrictamente una función de y con x fija, satisface todas las condiciones de una distribución de probabilidad. Esto también es cierto cuando $f(x, y)$ y $g(x)$ son la densidad conjunta y la distribución marginal, respectivamente, de variables aleatorias continuas. Como resultado, para poder calcular probabilidades condicionales de manera eficaz es sumamente importante que utilicemos el tipo especial de distribución de la forma $f(x, y)/g(x)$. Este tipo de distribución se llama **distribución de probabilidad condicional** y se define formalmente como sigue:

Definición 3.11: Sean X y Y dos variables aleatorias, discretas o continuas. La **distribución condicional** de la variable aleatoria Y , dado que $X = x$, es

$$f(y|x) = \frac{f(x, y)}{g(x)}, \text{ siempre que } g(x) > 0.$$

De manera similar, la distribución condicional de la variable aleatoria X , dado que $Y = y$, es

$$f(x|y) = \frac{f(x, y)}{h(y)}, \text{ siempre que } h(y) > 0.$$

Si deseamos encontrar la probabilidad de que la variable aleatoria discreta X caiga entre a y b cuando sabemos que la variable discreta $Y = y$, evaluamos

$$P(a < X < b | Y = y) = \sum_{a < x < b} f(x|y),$$

donde la sumatoria se extiende a todos los valores de X entre a y b . Cuando X y Y son continuas, evaluamos

$$P(a < X < b | Y = y) = \int_a^b f(x|y) dx.$$

Ejemplo 3.18: Remítase al ejemplo 3.14, calcule la distribución condicional de X , dado que $Y = 1$, y utilice el resultado para determinar $P(X = 0 | Y = 1)$.

Solución: Necesitamos encontrar $f(x|y)$, donde $y = 1$. Primero calculamos que

$$h(1) = \sum_{x=0}^2 f(x, 1) = \frac{3}{14} + \frac{3}{14} + 0 = \frac{3}{7}.$$

Ahora calculamos,

$$f(x|1) = \frac{f(x, 1)}{h(1)} = \left(\frac{7}{3}\right)f(x, 1), \quad x = 0, 1, 2.$$

Por lo tanto,

$$f(0|1) = \left(\frac{7}{3}\right)f(0, 1) = \left(\frac{7}{3}\right)\left(\frac{3}{14}\right) = \frac{1}{2}, \quad f(1|1) = \left(\frac{7}{3}\right)f(1, 1) = \left(\frac{7}{3}\right)\left(\frac{3}{14}\right) = \frac{1}{2},$$

$$f(2|1) = \left(\frac{7}{3}\right)f(2, 1) = \left(\frac{7}{3}\right)(0) = 0,$$

y la distribución condicional de X , dado que $Y = 1$, es

x	0	1	2
$f(x 1)$	$\frac{1}{2}$	$\frac{1}{2}$	0

Finalmente,

$$P(X = 0 | Y = 1) = f(0|1) = \frac{1}{2}.$$

Por lo tanto, si se sabe que 1 de los 2 repuestos seleccionados es rojo, tenemos una probabilidad igual a $1/2$ de que el otro repuesto no sea azul. \square

Ejemplo 3.19: La densidad conjunta para las variables aleatorias (X, Y) , donde X es el cambio unitario de temperatura y Y es la proporción de desplazamiento espectral que produce cierta partícula atómica es

$$f(x, y) = \begin{cases} 10xy^2, & 0 < x < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Calcule las densidades marginales $g(x)$, $h(y)$ y la densidad condicional $f(y|x)$.
 b) Calcule la probabilidad de que el espectro se desplace más de la mitad de las observaciones totales, dado que la temperatura se incremente en 0.25 unidades.

Solución: a) Por definición,

$$g(x) = \int_{-\infty}^{\infty} f(x, y) dy = \int_x^1 10xy^2 dy$$

$$= \frac{10}{3}xy^3 \Big|_{y=x}^{y=1} = \frac{10}{3}x(1 - x^3), \quad 0 < x < 1,$$

$$h(y) = \int_{-\infty}^{\infty} f(x, y) dx = \int_0^y 10xy^2 dx = 5x^2y^2 \Big|_{x=0}^{x=y} = 5y^4, \quad 0 < y < 1.$$

Entonces,

$$f(y|x) = \frac{f(x, y)}{g(x)} = \frac{10xy^2}{\frac{10}{3}x(1 - x^3)} = \frac{3y^2}{1 - x^3}, \quad 0 < x < y < 1.$$

b) Por lo tanto,

$$P\left(Y > \frac{1}{2} \mid X = 0.25\right) = \int_{1/2}^1 f(y | x = 0.25) dy = \int_{1/2}^1 \frac{3y^2}{1 - 0.25^3} dy = \frac{8}{9}. \quad \square$$

Ejemplo 3.20: Dada la función de densidad conjunta

$$f(x, y) = \begin{cases} \frac{x(1+3y^2)}{4}, & 0 < x < 2, \quad 0 < y < 1, \\ 0, & \text{en otro caso,} \end{cases}$$

calcule $g(x)$, $h(y)$, $f(x|y)$ y evalúe $P\left(\frac{1}{4} < X < \frac{1}{2} \mid Y = \frac{1}{3}\right)$.

Solución: Por definición de la densidad marginal, para $0 < x < 2$,

$$\begin{aligned} g(x) &= \int_{-\infty}^{\infty} f(x, y) dy = \int_0^1 \frac{x(1+3y^2)}{4} dy \\ &= \left(\frac{xy}{4} + \frac{xy^3}{4} \right) \Big|_{y=0}^{y=1} = \frac{x}{2}, \end{aligned}$$

y para $0 < y < 1$,

$$\begin{aligned} h(y) &= \int_{-\infty}^{\infty} f(x, y) dx = \int_0^2 \frac{x(1+3y^2)}{4} dx \\ &= \left(\frac{x^2}{8} + \frac{3x^2y^2}{8} \right) \Big|_{x=0}^{x=2} = \frac{1+3y^2}{2}. \end{aligned}$$

Por lo tanto, usando la definición de la densidad condicional para $0 < x < 2$,

$$f(x|y) = \frac{f(x, y)}{h(y)} = \frac{x(1+3y^2)/4}{(1+3y^2)/2} = \frac{x}{2},$$

y

$$P\left(\frac{1}{4} < X < \frac{1}{2} \mid Y = \frac{1}{3}\right) = \int_{1/4}^{1/2} \frac{x}{2} dx = \frac{3}{64}.$$

Independencia estadística

Si $f(x|y)$ no depende de y , como ocurre en el ejemplo 3.20, entonces $f(x|y) = g(x)$ y $f(x, y) = g(x)h(y)$. La prueba se realiza sustituyendo

$$f(x, y) = f(x|y)h(y)$$

en la distribución marginal de X . Es decir,

$$g(x) = \int_{-\infty}^{\infty} f(x, y) dy = \int_{-\infty}^{\infty} f(x|y)h(y) dy.$$

Si $f(x|y)$ no depende de y , podemos escribir

$$g(x) = f(x|y) \int_{-\infty}^{\infty} h(y) dy.$$

Entonces,

$$\int_{-\infty}^{\infty} h(y) dy = 1,$$

ya que $h(y)$ es la función de densidad de probabilidad de Y . Por lo tanto,

$$g(x) = f(x|y) \text{ y entonces } f(x, y) = g(x)h(y).$$

Debería tener sentido para el lector que si $f(x|y)$ no depende de y , entonces, por supuesto, el resultado de la variable aleatoria Y no repercute en el resultado de la variable aleatoria X . En otras palabras, decimos que X y Y son variables aleatorias independientes. Ofrecemos ahora la siguiente definición formal de independencia estadística.

Definición 3.12: Sean X y Y dos variables aleatorias, discretas o continuas, con distribución de probabilidad conjunta $f(x,y)$ y distribuciones marginales $g(x)$ y $h(y)$, respectivamente. Se dice que las variables aleatorias X y Y son **estadísticamente independientes** si y sólo si

$$f(x, y) = g(x)h(y)$$

para toda (x,y) dentro de sus rangos.

Las variables aleatorias continuas del ejemplo 3.20 son estadísticamente independientes, pues el producto de las dos distribuciones marginales da la función de densidad conjunta. Sin embargo, es evidente que ése no es el caso de las variables continuas del ejemplo 3.19. La comprobación de la independencia estadística de variables aleatorias discretas requiere una investigación más profunda, ya que es posible que el producto de las distribuciones marginales sea igual a la distribución de probabilidad conjunta para algunas, aunque no para todas, las combinaciones de (x,y) . Si puede encontrar algún punto (x,y) para el que $f(x,y)$ se define de manera que $f(x,y) \neq g(x)h(y)$, las variables discretas X y Y no son estadísticamente independientes.

Ejemplo 3.21: Demuestre que las variables aleatorias del ejemplo 3.14 no son estadísticamente independientes.

Prueba: Consideremos el punto $(0,1)$. A partir de la tabla 3.1, encontramos que las tres probabilidades $f(0,1)$, $g(0)$ y $h(1)$ son

$$\begin{aligned} f(0, 1) &= \frac{3}{14}, \\ g(0) &= \sum_{y=0}^2 f(0, y) = \frac{3}{28} + \frac{3}{14} + \frac{1}{28} = \frac{5}{14}, \\ h(1) &= \sum_{x=0}^2 f(x, 1) = \frac{3}{14} + \frac{3}{14} + 0 = \frac{3}{7}. \end{aligned}$$

Claramente,

$$f(0, 1) \neq g(0)h(1),$$

por lo tanto, X y Y no son estadísticamente independientes. ▮

Todas las definiciones anteriores respecto a dos variables aleatorias se pueden generalizar al caso de n variables aleatorias. Sea $f(x_1, x_2, \dots, x_n)$ la función de probabilidad conjunta de las variables aleatorias X_1, X_2, \dots, X_n . La distribución marginal de X_1 , por ejemplo, es

$$g(x_1) = \sum_{x_2} \cdots \sum_{x_n} f(x_1, x_2, \dots, x_n)$$

para el caso discreto, y

$$g(x_1) = \sum_{x_2} \cdots \sum_{x_n} f(x_1, x_2, \dots, x_n)$$

para el caso continuo. Ahora podemos obtener **distribuciones marginales conjuntas** como $g(x_1, x_2)$, donde

$$g(x_1, x_2) = \begin{cases} \sum_{x_3} \cdots \sum_{x_n} f(x_1, x_2, \dots, x_n) & \text{(caso discreto),} \\ \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, x_2, \dots, x_n) dx_3 dx_4 \cdots dx_n & \text{(caso continuo).} \end{cases}$$

Podríamos considerar numerosas distribuciones condicionales. Por ejemplo, la **distribución condicional conjunta** de X_1, X_2 y X_3 , dado que $X_4 = x_4, X_5 = x_5, \dots, X_n = x_n$, se escribe como

$$f(x_1, x_2, x_3 \mid x_4, x_5, \dots, x_n) = \frac{f(x_1, x_2, \dots, x_n)}{g(x_4, x_5, \dots, x_n)},$$

donde $g(x_4, x_5, \dots, x_n)$ es la distribución marginal conjunta de las variables aleatorias X_4, X_5, \dots, X_n .

Una generalización de la definición 3.12 nos lleva a la siguiente definición para la independencia estadística mutua de las variables X_1, X_2, \dots, X_n .

Definición 3.13: Sean X_1, X_2, \dots, X_n n variables aleatorias, discretas o continuas, con distribución de probabilidad conjunta $f(x_1, x_2, \dots, x_n)$ y distribuciones marginales $f_1(x_1), f_2(x_2), \dots, f_n(x_n)$, respectivamente. Se dice que las variables aleatorias X_1, X_2, \dots, X_n son **recíproca y estadísticamente independientes** si y sólo si

$$f(x_1, x_2, \dots, x_n) = f_1(x_1)f_2(x_2) \cdots f_n(x_n)$$

para toda (x_1, x_2, \dots, x_n) dentro de sus rangos.

Ejemplo 3.22: Suponga que el tiempo de vida en anaquel de cierto producto comestible precedero empacado en cajas de cartón, en años, es una variable aleatoria cuya función de densidad de probabilidad está dada por

$$f(x) = \begin{cases} e^{-x} & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

Represente los tiempos de vida en anaquel para tres de estas cajas seleccionadas de forma independiente con X_1, X_2 y X_3 y calcule $P(X_1 < 2, 1 < X_2 < 3, X_3 > 2)$.

Solución: Como las cajas se seleccionaron de forma independiente, suponemos que las variables aleatorias X_1, X_2 y X_3 son estadísticamente independientes y que tienen la siguiente densidad de probabilidad conjunta:

$$f(x_1, x_2, x_3) = f(x_1)f(x_2)f(x_3) = e^{-x_1}e^{-x_2}e^{-x_3} = e^{-x_1-x_2-x_3},$$

para $x_1 > 0, x_2 > 0, x_3 > 0$, y $f(x_1, x_2, x_3) = 0$ en cualquier otro caso. En consecuencia,

$$\begin{aligned} P(X_1 < 2, 1 < X_2 < 3, X_3 > 2) &= \int_2^{\infty} \int_1^3 \int_0^2 e^{-x_1-x_2-x_3} dx_1 dx_2 dx_3 \\ &= (1 - e^{-2})(e^{-1} - e^{-3})e^{-2} = 0.0372. \end{aligned}$$

¿Por qué son importantes las características de las distribuciones de probabilidad y de dónde provienen?

Es importante que este texto ofrezca al lector una transición hacia los siguientes tres capítulos. En los ejemplos y los ejercicios presentamos casos de situaciones prácticas de ingeniería y ciencias, en los cuales las distribuciones de probabilidad y sus propiedades se utilizan para resolver problemas importantes. Estas distribuciones de probabilidad, ya sean discretas o continuas, se presentaron mediante frases como “se sabe que”, “suponga que” o incluso, en ciertos casos, “la evidencia histórica sugiere que”. Se trata de situaciones en las que la naturaleza de la distribución, e incluso una estimación óptima de la estructura de la probabilidad, se pueden determinar utilizando datos históricos, datos tomados de estudios a largo plazo o incluso de grandes cantidades de datos planeados. El lector debería tener presente lo expuesto en el capítulo 1 respecto al uso de histogramas y, por consiguiente, recordar cómo se estiman las distribuciones de frecuencias a partir de los histogramas. Sin embargo, no todas las funciones de probabilidad y de densidad de probabilidad se derivan de cantidades grandes de datos históricos. Hay un gran número de situaciones en las que la naturaleza del escenario científico sugiere un tipo de distribución. De hecho, varias de ellas se reflejan en los ejercicios del capítulo 2 y en este capítulo. Cuando observaciones repetidas independientes son binarias por naturaleza (es decir, defectuoso o no, funciona o no, alérgico o no) con un valor de 0 o 1, la distribución que cubre esta situación se llama **distribución binomial**. La función de probabilidad de esta distribución se explicará y se demostrará en el capítulo 5. El ejercicio 3.34 de la sección 3.3 y el ejercicio de repaso 3.80 constituyen ejemplos de este tipo de distribución, y hay otros que el lector también debería reconocer. El escenario de una distribución continua del tiempo de operación antes de cualquier falla, como en el ejercicio de repaso 3.69 o en el ejercicio 3.27 de la página 93, a menudo sugiere una clase de distribución denominada **distribución exponencial**. Tales tipos de ejemplos son tan sólo dos de la gran cantidad de las llamadas distribuciones estándar que se utilizan ampliamente en situaciones del mundo real porque el escenario científico que da lugar a cada uno de ellos es reconocible y a menudo se presenta en la práctica. Los capítulos 5 y 6 abarcan muchos de estos tipos de ejemplos, junto con alguna teoría inherente respecto de su uso.

La segunda parte de esta transición al material de los capítulos siguientes tiene que ver con el concepto de **parámetros de la población** o **parámetros de distribución**. Recuerde que en el capítulo 1 analizamos la necesidad de utilizar datos para ofrecer información sobre dichos parámetros. Profundizamos en el estudio de las nociones de **media** y de **varianza**, y proporcionamos ideas sobre esos conceptos en el contexto de una población. De hecho, es fácil calcular la media y la varianza de la población a partir de la función de probabilidad para el caso discreto, o de la función de densidad de probabilidad para el caso continuo. Tales parámetros y su importancia en la solución de muchas clases de problemas de la vida real nos proporcionarán gran parte del material de los capítulos 8 a 17.

Ejercicios

3.37 Determine los valores de c , tales que las siguientes funciones representen distribuciones de probabilidad conjunta de las variables aleatorias X y Y :

a) $f(x, y) = cxy$, para $x = 1, 2, 3$; $y = 1, 2, 3$;

b) $f(x, y) = c|x - y|$, para $x = -2, 0, 2$; $y = -2, 3$.

3.38 Si la distribución de probabilidad conjunta de X y Y está dada por

$$f(x, y) = \frac{x + y}{30}, \quad \text{para } x = 0, 1, 2, 3; y = 0, 1, 2.$$

calcule

- a) $P(X \leq 2, Y = 1)$;
 b) $P(X > 2, Y \leq 1)$;
 c) $P(X > Y)$;
 d) $P(X + Y = 4)$.

3.39 De un saco de frutas que contiene 3 naranjas, 2 manzanas y 3 plátanos se selecciona una muestra aleatoria de 4 frutas. Si X es el número de naranjas y Y el de manzanas en la muestra, calcule

- a) la distribución de probabilidad conjunta de X y Y ;
 b) $P\{(X, Y) \in A\}$, donde A es la región dada por $\{(x, y) | x + y \leq 2\}$.

3.40 Un restaurante de comida rápida opera tanto en un local que da servicio en el automóvil, como en un local que atiende a los clientes que llegan caminando. En un día elegido al azar, represente las proporciones de tiempo que el primero y el segundo local están en servicio con X y Y , respectivamente, y suponga que la función de densidad conjunta de estas variables aleatorias es

$$f(x, y) = \begin{cases} \frac{2}{3}(x + 2y), & 0 \leq x \leq 1, 0 \leq y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Calcule la densidad marginal de X .
 b) Calcule la densidad marginal de Y .
 c) Calcule la probabilidad de que el local que da servicio a los clientes que llegan en automóvil esté lleno menos de la mitad del tiempo.

3.41 Una empresa dulcera distribuye cajas de chocolates con un surtido de cremas, chiclosos y envinados. Suponga que cada caja pesa 1 kilogramo, pero que los pesos individuales de cremas, chiclosos y envinados varían de una a otra cajas. Para una caja seleccionada al azar, represente los pesos de las cremas y los chiclosos con X y Y , respectivamente, y suponga que la función de densidad conjunta de estas variables es

$$f(x, y) = \begin{cases} 24xy, & 0 \leq x \leq 1, 0 \leq y \leq 1, x + y \leq 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

- a) Calcule la probabilidad de que en una caja dada los envinados representen más de la mitad del peso.
 b) Calcule la densidad marginal para el peso de las cremas.
 c) Calcule la probabilidad de que el peso de los chiclosos en una caja sea menor que $1/8$ de kilogramo, si se sabe que las cremas constituyen $3/4$ partes del peso.

3.42 Sean X y Y la duración de la vida, en años, de dos componentes en un sistema electrónico. Si la función de densidad conjunta de estas variables es

$$f(x, y) = \begin{cases} e^{-(x+y)}, & x > 0, y > 0, \\ 0, & \text{en otro caso,} \end{cases}$$

calcule $P(0 < X < 1 | Y = 2)$.

3.43 Sea X el tiempo de reacción, en segundos, ante cierto estímulo, y Y la temperatura (en °F) a la cual inicia cierta reacción. Suponga que dos variables aleatorias, X y Y , tienen la densidad conjunta

$$f(x, y) = \begin{cases} 4xy, & 0 < x < 1, 0 < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule

- a) $P(0 \leq X \leq \frac{1}{2} \text{ y } \frac{1}{4} \leq Y \leq \frac{1}{2})$;
 b) $P(X < Y)$.

3.44 Se supone que cada rueda trasera de un avión experimental se llena a una presión de 40 libras por pulgada cuadrada (psi). Sea X la presión real del aire para la rueda derecha y Y la presión real del aire de la rueda izquierda. Suponga que X y Y son variables aleatorias con la siguiente función de densidad conjunta:

$$f(x, y) = \begin{cases} k(x^2 + y^2), & 30 \leq x < 50, 30 \leq y < 50, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Calcule k .
 b) Calcule $P(30 \leq X \leq 40 \text{ y } 40 \leq Y < 50)$.
 c) Calcule la probabilidad de que ambas ruedas no contengan la suficiente cantidad de aire.

3.45 Sea X el diámetro de un cable eléctrico blindado y Y el diámetro del molde cerámico que hace el cable. Tanto X como Y tienen una escala tal que están entre 0 y 1. Suponga que X y Y tienen la siguiente densidad conjunta:

$$f(x, y) = \begin{cases} \frac{1}{y}, & 0 < x < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule $P(X + Y > 1/2)$.

3.46 Remítase al ejercicio 3.38, calcule

- a) la distribución marginal de X ;
 b) la distribución marginal de Y .

3.47 Al principio de cualquier día la cantidad de que-roso que contiene un tanque, en miles de litros, es una cantidad aleatoria Y , de la que durante el día se vende una cantidad aleatoria X . Suponga que el tanque no se reabastece durante el día, de manera que $x \leq y$, e imagine también que la función de densidad conjunta de estas variables es

$$f(x, y) = \begin{cases} 2, & 0 < x \leq y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Determine si X y Y son independientes.
 b) Calcule $P(1/4 < X < 1/2 | Y = 3/4)$.

3.48 Remítase al ejercicio 3.39 y calcule

- a) $f(y|2)$ para todos los valores de y ;
 b) $P(Y=0|X=2)$.

3.49 Sea X el número de veces que fallará cierta máquina de control numérico: 1, 2 o 3 veces en un día dado. Y si Y denota el número de veces que se llama a un técnico para una emergencia, su distribución de probabilidad conjunta estará dada como

$f(x, y)$	x		
	1	2	3
1	0.05	0.05	0.10
y	3	0.05	0.10
5	0.00	0.20	0.10

- a) Evalúe la distribución marginal de X .
 b) Evalúe la distribución marginal de Y .
 c) Calcule $P(Y=3|X=2)$.

3.50 Suponga que X y Y tienen la siguiente distribución de probabilidad conjunta:

$f(x, y)$	x	
	2	4
1	0.10	0.15
y	3	0.20
5	0.10	0.15

- a) Calcule la distribución marginal de X .
 b) Calcule la distribución marginal de Y .

3.51 De las 12 cartas mayores (jotas, reinas y reyes) de una baraja ordinaria de 52 cartas se sacan tres cartas sin reemplazo. Sea X el número de reyes que se seleccionan y Y el número de jotas. Calcule

- a) la distribución de probabilidad conjunta de X y Y ;
 b) $P[(X, Y) \in A]$, donde A es la región dada por $\{(x, y) | x + y \geq 2\}$.

3.52 Una moneda se lanza dos veces. Sea Z el número de caras en el primer lanzamiento y W el número total de caras en los 2 lanzamientos. Si la moneda no está balanceada y una cara tiene una probabilidad de ocurrencia de 40%, calcule

- a) la distribución de probabilidad conjunta de W y Z ;
 b) la distribución marginal de W ;
 c) la distribución marginal de Z ;
 d) la probabilidad de que ocurra al menos 1 cara.

3.53 Dada la función de densidad conjunta

$$f(x, y) = \begin{cases} \frac{6-x-y}{8}, & 0 < x < 2, 2 < y < 4, \\ 0, & \text{en otro caso,} \end{cases}$$

calcule $P(1 < Y < 3 | X = 1)$.

3.54 Determine si las dos variables aleatorias del ejercicio 3.49 son dependientes o independientes.

3.55 Determine si las dos variables aleatorias del ejercicio 3.50 son dependientes o independientes.

3.56 La función de densidad conjunta de las variables aleatorias X y Y es

$$f(x, y) = \begin{cases} 6x, & 0 < x < 1, 0 < y < 1-x, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Demuestre que X y Y no son independientes.
 b) Calcule $P(X > 0.3 | Y = 0.5)$.

3.57 Si X , Y y Z tienen la siguiente función de densidad de probabilidad conjunta:

$$f(x, y, z) = \begin{cases} kxy^2z, & 0 < x, y < 1, 0 < z < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Calcule k .
 b) Calcule $P(X < \frac{1}{4}, Y > \frac{1}{2}, 1 < Z < 2)$.

3.58 Determine si las dos variables aleatorias del ejercicio 3.43 son dependientes o independientes.

3.59 Determine si las dos variables aleatorias del ejercicio 3.44 son dependientes o independientes.

3.60 La función de densidad de probabilidad conjunta de las variables aleatorias X , Y y Z es

$$f(x, y, z) = \begin{cases} \frac{4xyz^2}{9}, & 0 < x, y < 1, 0 < z < 3, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule

- a) la función de densidad marginal conjunta de Y y Z ;
 b) la densidad marginal de Y ;
 c) $P(\frac{1}{4} < X < \frac{1}{2}, Y > \frac{1}{4}, 1 < Z < 2)$;
 d) $P(0 < X < \frac{1}{2} | Y = \frac{1}{4}, Z = 2)$.

Ejercicios de repaso

3.61 Una empresa tabacalera produce mezclas de tabaco. Cada mezcla contiene diversas proporciones de tabaco turco, tabaco de la región y otros. Las proporciones de tabaco turco y de la región en una mezcla son variables aleatorias con una función de densidad conjunta ($X = \text{turco}$ y $Y = \text{de la región}$)

$$f(x, y) = \begin{cases} 24xy, & 0 \leq x, y \leq 1, x + y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Calcule la probabilidad de que en determinada caja el tabaco turco represente más de la mitad de la mezcla.
- Calcule la función de densidad marginal para la proporción del tabaco de la región.
- Calcule la probabilidad de que la proporción de tabaco turco sea menor que $1/8$, si se sabe que la mezcla contiene $3/4$ de tabaco de la región.

3.62 Una empresa de seguros ofrece a sus asegurados varias opciones diferentes de pago de la prima. Para un asegurado seleccionado al azar, sea X el número de meses entre pagos sucesivos. La función de distribución acumulada de X es

$$F(x) = \begin{cases} 0, & \text{si } x < 1, \\ 0.4, & \text{si } 1 \leq x < 3, \\ 0.6, & \text{si } 3 \leq x < 5, \\ 0.8, & \text{si } 5 \leq x < 7, \\ 1.0, & \text{si } x \geq 7. \end{cases}$$

- ¿Cuál es la función de masa de probabilidad de X ?
- Calcule $P(4 < X \leq 7)$.

3.63 Dos componentes electrónicos de un sistema de proyectiles funcionan en conjunto para el éxito de todo el sistema. Sean X y Y la vida en horas de los dos componentes. La densidad conjunta de X y Y es

$$f(x, y) = \begin{cases} ye^{-y(1+x)}, & x, y \geq 0, \\ 0, & \text{en otro caso.} \end{cases}$$

- Determine las funciones de densidad marginal para ambas variables aleatorias.
- ¿Cuál es la probabilidad de que ambos componentes duren más de dos horas?

3.64 Una instalación de servicio opera con dos líneas telefónicas. En un día elegido al azar, sea X la proporción de tiempo que la primera línea está en uso, mientras que Y es la proporción de tiempo en que la segunda línea está en uso. Suponga que la función de densidad de probabilidad conjunta para (X, Y) es

$$f(x, y) = \begin{cases} \frac{3}{2}(x^2 + y^2), & 0 \leq x, y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Calcule la probabilidad de que ninguna línea esté ocupada más de la mitad del tiempo.
- Calcule la probabilidad de que la primera línea esté ocupada más del 75% del tiempo.

3.65 Sea el número de llamadas telefónicas que recibe un conmutador durante un intervalo de 5 minutos una variable aleatoria X con la siguiente función de probabilidad:

$$f(x) = \frac{e^{-2}2^x}{x!}, \quad \text{para } x = 0, 1, 2, \dots$$

- Determine la probabilidad de que X sea igual a 0, 1, 2, 3, 4, 5 y 6.
- Grafique la función de masa de probabilidad para estos valores de x .
- Determine la función de distribución acumulada para estos valores de X .

3.66 Considere las variables aleatorias X y Y con la siguiente función de densidad conjunta

$$f(x, y) = \begin{cases} x + y, & 0 \leq x, y \leq 1, \\ 0, & \text{en cualquier otro caso.} \end{cases}$$

- Calcule las distribuciones marginales de X y Y .
- Calcule $P(X > 0.5, Y > 0.5)$.

3.67 En un proceso industrial se elaboran artículos que se pueden clasificar como defectuosos o no defectuosos. La probabilidad de que un artículo esté defectuoso es 0.1. Se realiza un experimento en el que 5 artículos del proceso se eligen al azar. Sea la variable aleatoria X el número de artículos defectuosos en esta muestra de 5. ¿Cuál es la función de masa de probabilidad de X ?

3.68 Considere la siguiente función de densidad de probabilidad conjunta de las variables aleatorias X y Y :

$$f(x, y) = \begin{cases} \frac{3x-y}{9}, & 1 < x < 3, 1 < y < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

- Calcule las funciones de densidad marginal de X y Y .
- ¿ X y Y son independientes?
- Calcule $P(X > 2)$.

3.69 La duración en horas de un componente eléctrico es una variable aleatoria con la siguiente función de distribución acumulada:

$$F(x) = \begin{cases} 1 - e^{-\frac{x}{50}}, & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Determine su función de densidad de probabilidad.
 b) Determine la probabilidad de que la vida útil de tal componente rebase las 70 horas.

3.70 En una fábrica específica de pantalones un grupo de 10 trabajadores los inspecciona tomando aleatoriamente algunos de la línea de producción. A cada inspector se le asigna un número del 1 al 10. Un comprador selecciona un pantalón para adquirirlo. Sea la variable aleatoria X el número del inspector.

- a) Determine una función de masa de probabilidad razonable para X .
 b) Grafique la función de distribución acumulada para X .

3.71 La vida en anaquel de un producto es una variable aleatoria que se relaciona con la aceptación por parte del consumidor. Resulta que la vida en anaquel Y , en días, de cierta clase de artículo de panadería tiene la siguiente función de densidad:

$$f(y) = \begin{cases} \frac{1}{2}e^{-y/2}, & 0 \leq y < \infty, \\ 0, & \text{en otro caso.} \end{cases}$$

¿Qué fracción de las rebanadas de este producto que hoy están en exhibición se espera que se vendan en 3 días a partir de hoy?

3.72 El congestionamiento de pasajeros es un problema de servicio en los aeropuertos, en los cuales se instalan trenes para reducir la congestión. Cuando se usa el tren, el tiempo X , en minutos, que toma viajar desde la terminal principal hasta una explanada específica tiene la siguiente función de densidad:

$$f(x) = \begin{cases} \frac{1}{10}, & 0 \leq x \leq 10, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Demuestre que la función de densidad de probabilidad anterior es válida.
 b) Calcule la probabilidad de que el tiempo que le toma a un pasajero viajar desde la terminal principal hasta la explanada no exceda los 7 minutos.

3.73 Las impurezas en el lote del producto final de un proceso químico a menudo reflejan un grave problema. A partir de una cantidad considerable de datos recabados en la planta se sabe que la proporción Y de impurezas en un lote tiene una función de densidad dada por

$$f(y) = \begin{cases} 10(1-y)^9, & 0 \leq y \leq 1, \\ 0, & \text{en cualquier otro caso.} \end{cases}$$

- a) Verifique que la función de densidad anterior sea válida.
 b) Se considera que un lote no es vendible y, por consiguiente, no es aceptable si el porcentaje de impurezas es superior a 60%. Con la calidad del proceso

actual, ¿cuál es el porcentaje de lotes que no son aceptables?

3.74 El tiempo Z , en minutos, entre llamadas a un sistema de alimentación eléctrica tiene la siguiente función de densidad de probabilidad:

$$f(z) = \begin{cases} \frac{1}{10}e^{-z/10}, & 0 < z < \infty, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) ¿Cuál es la probabilidad de que no haya llamadas en un lapso de 20 minutos?
 b) ¿Cuál es la probabilidad de que la primera llamada entre en los primeros 10 minutos después de abrir?

3.75 Un sistema químico que surge de una reacción química tiene dos componentes importantes, entre otros, en una mezcla. La distribución conjunta que describe las proporciones X_1 y X_2 de estos dos componentes está dada por

$$f(x_1, x_2) = \begin{cases} 2, & 0 < x_1 < x_2 < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Determine la distribución marginal de X_1 .
 b) Determine la distribución marginal de X_2 .
 c) ¿Cuál es la probabilidad de que las proporciones del componente generen los resultados $X_1 < 0.2$ y $X_2 > 0.5$?
 d) Determine la distribución condicional $f_{X_1|X_2}(x_1 | x_2)$.

3.76 Considere la situación del ejercicio de repaso 3.75; pero suponga que la distribución conjunta de las dos proporciones está dada por

$$f(x_1, x_2) = \begin{cases} 6x_2, & 0 < x_2 < x_1 < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Determine la distribución marginal $f_{X_1}(x_1)$ de la proporción X_1 y verifique que sea una función de densidad válida.
 b) ¿Cuál es la probabilidad de que la proporción X_2 sea menor que 0.5, dado que X_1 es 0.7?

3.77 Considere las variables aleatorias X y Y que representan el número de vehículos que llegan a dos esquinas de calles separadas durante cierto periodo de 2 minutos. Estas esquinas de las calles están bastante cerca una de la otra, así que es importante que los ingenieros de tráfico se ocupen de ellas de manera conjunta si fuera necesario. Se sabe que la distribución conjunta de X y Y es

$$f(x, y) = \frac{9}{16} \cdot \frac{1}{4^{(x+y)}},$$

para $x = 0, 1, 2, \dots$, y para $y = 0, 1, 2, \dots$

- a) ¿Son independientes las dos variables aleatorias X y Y ? Explique su respuesta.

b) ¿Cuál es la probabilidad de que, durante el periodo en cuestión, lleguen menos de 4 vehículos a las dos esquinas?

3.78 El comportamiento de series de componentes desempeña un papel importante en problemas de confiabilidad científicos y de ingeniería. Ciertamente la confiabilidad de todo el sistema no es mejor que la del componente más débil de las series. En un sistema de series los componentes funcionan de manera independiente unos de otros. En un sistema particular de tres componentes, la probabilidad de cumplir con la especificación para los componentes 1, 2 y 3, respectivamente, son 0.95, 0.99 y 0.92. ¿Cuál es la probabilidad de que todo el sistema funcione?

3.79 Otro tipo de sistema que se utiliza en trabajos de ingeniería es un grupo de componentes en paralelo o un sistema paralelo. En este enfoque más conservador la probabilidad de que el sistema funcione es mayor que la probabilidad de que cualquier componente funcione. El sistema fallará sólo cuando falle todo el sistema. Considere una situación en la que hay 4 componentes

independientes en un sistema paralelo, en la que la probabilidad de operación está dada por

Componente 1: 0.95; Componente 2: 0.94;
Componente 3: 0.90; Componente 4: 0.97.

¿Cuál es la probabilidad de que no falle el sistema?

3.80 Considere un sistema de componentes en que hay cinco componentes independientes, cada uno de los cuales tiene una probabilidad de operación de 0.92. De hecho, el sistema tiene una redundancia preventiva diseñada para que no falle mientras 3 de sus 5 componentes estén en funcionamiento. ¿Cuál es la probabilidad de que funcione todo el sistema?

3.81 Proyecto de grupo: Observe el color de los zapatos de los estudiantes en 5 periodos de clases. Suponga que las categorías de color son rojo, blanco, negro, café y otro. Construya una tabla de frecuencias para cada color.

- Estime e interprete el significado de la distribución de probabilidad.
- ¿Cuál es la probabilidad estimada de que en el siguiente periodo de clases un estudiante elegido al azar use un par de zapatos rojos o blancos?

3.5 Posibles riesgos y errores conceptuales; relación con el material de otros capítulos

En los siguientes capítulos será evidente que las distribuciones de probabilidad representan la estructura mediante la cual las probabilidades que se calculan ayudan a evaluar y a comprender un proceso. Por ejemplo, en el ejercicio de repaso 3.65 la distribución de probabilidad que cuantifica la probabilidad de que haya una carga excesiva durante ciertos periodos podría ser muy útil en la planeación de cualquier cambio en el sistema. El ejercicio de repaso 3.69 describe un escenario donde se estudia el periodo de vida útil de un componente electrónico. Conocer la estructura de la probabilidad para el componente contribuirá de manera significativa al entendimiento de la confiabilidad de un sistema mayor del cual éste forme parte. Además, comprender la naturaleza general de las distribuciones de probabilidad reforzará el conocimiento del concepto **valor- P** , que se estudió brevemente en el capítulo 1 y que desempeñará un papel destacado al inicio del capítulo 10 y en lo que resta del texto.

Los capítulos 4, 5 y 6 dependen mucho del material cubierto en este capítulo. En el capítulo 4 estudiaremos el significado de **parámetros** importantes en las distribuciones de probabilidad. Tales parámetros cuantifican las nociones de **tendencia central** y **variabilidad** en un sistema. De hecho, el conocimiento de tales cantidades, al margen de la distribución completa, puede ofrecer información sobre la naturaleza del sistema. En los capítulos 5 y 6 se examinarán escenarios de ingeniería, biológicos y de ciencia en general que identifican tipos de distribuciones especiales. Por ejemplo, la estructura de la función de probabilidad en el ejercicio de repaso 3.65 se identificará fácilmente bajo ciertas suposiciones que se estudiarán en el capítulo 5. Lo mismo ocurre en el contexto

del ejercicio de repaso 3.69, que es un caso especial de problema sobre **tiempo de operación antes de la falla**, cuya función de densidad de probabilidad se estudiará en el capítulo 6.

En lo que concierne a los riesgos potenciales de utilizar el material de este capítulo, la advertencia para el lector sería no interpretar el material más allá de lo que sea evidente. La naturaleza general de la distribución de probabilidad para un fenómeno científico determinado no es obvia a partir de lo que se estudió aquí. La finalidad de este capítulo es que los lectores aprendan a manipular una distribución de probabilidad, no que aprendan a identificar un tipo específico. Los capítulos 5 y 6 avanzan un largo trecho hacia la identificación de acuerdo con la naturaleza general del sistema científico.

Capítulo 4

Esperanza matemática

4.1 Media de una variable aleatoria

En el capítulo 1 estudiamos la media muestral, que es la media aritmética de los datos. Ahora considere la siguiente situación: si dos monedas se lanzan 16 veces y X es el número de caras que resultan en cada lanzamiento, entonces los valores de X pueden ser 0, 1 y 2. Suponga que los resultados del experimento son: cero caras, una cara y dos caras, un total de 4, 7 y 5 veces, respectivamente. El número promedio de caras por lanzamiento de las dos monedas es, entonces,

$$\frac{(0)(4) + (1)(7) + (2)(5)}{16} = 1.06.$$

Éste es un valor promedio de los datos, aunque no es un resultado posible de $\{0, 1, 2\}$. Por lo tanto, un promedio no es necesariamente un resultado posible del experimento. Por ejemplo, es probable que el ingreso mensual promedio de un vendedor no sea igual a alguno de sus cheques de pago mensuales.

Reestructuremos ahora nuestro cálculo del número promedio de caras para tener la siguiente forma equivalente:

$$(0) \left(\frac{4}{16} \right) + (1) \left(\frac{7}{16} \right) + (2) \left(\frac{5}{16} \right) = 1.06.$$

Los números $4/16$, $7/16$ y $5/16$ son las fracciones de los lanzamientos totales que dan como resultado 0, 1 y 2 caras, respectivamente. Tales fracciones también son las frecuencias relativas de los diferentes valores de X en nuestro experimento. Entonces, realmente podemos calcular la media, o el promedio de un conjunto de datos, si conocemos los distintos valores que ocurren y sus frecuencias relativas sin tener conocimiento del número total de observaciones en el conjunto de datos. Por lo tanto, si $4/16$ o $1/4$ de los lanzamientos dan como resultado cero caras, $7/16$ de los lanzamientos dan como resultado una cara y $5/16$ dan como resultado dos caras, el número medio de caras por lanzamiento sería 1.06, sin importar si el número total de lanzamientos fue 16, 1000 o incluso 10,000.

Este método de frecuencias relativas se utiliza para calcular el número promedio de caras que esperaríamos obtener a largo plazo por el lanzamiento de dos monedas. A este valor promedio se le conoce como **media de la variable aleatoria X** o **media de la distribución de probabilidad de X** , y se le denota como μ_x o simplemente como μ cuando es evidente a qué variable aleatoria se está haciendo referencia. También es común entre los estadísticos referirse a esta media como la **esperanza matemática** o el **valor esperado** de la variable aleatoria X y denotarla como $E(X)$.

Suponiendo que una moneda legal se lanza dos veces, encontramos que el espacio muestral para el experimento es

$$S = \{HH, HT, TH, TT\}.$$

Como los 4 puntos muestrales son igualmente probables, se deduce que

$$P(X = 0) = P(TT) = \frac{1}{4}, \quad P(X = 1) = P(TH) + P(HT) = \frac{1}{2},$$

y

$$P(X = 2) = P(HH) = \frac{1}{4},$$

donde un elemento típico, digamos TH , indica que el primer lanzamiento dio como resultado una cruz seguida por una cara en el segundo lanzamiento. Así, estas probabilidades son precisamente las frecuencias relativas para los eventos dados a largo plazo. Por lo tanto,

$$\mu = E(X) = (0) \left(\frac{1}{4}\right) + (1) \left(\frac{1}{2}\right) + (2) \left(\frac{1}{4}\right) = 1.$$

Este resultado significa que una persona que lance 2 monedas una y otra vez obtendrá, en promedio, 1 cara por cada lanzamiento.

El método descrito antes para calcular el número esperado de caras cada vez que se lanzan 2 monedas sugiere que la media, o el valor esperado de cualquier variable aleatoria discreta, se puede obtener multiplicando cada uno de los valores x_1, x_2, \dots, x_n de la variable aleatoria X por su probabilidad correspondiente $f(x_1), f(x_2), \dots, f(x_n)$ y sumando los productos. Esto es cierto, sin embargo, sólo si la variable aleatoria es discreta. En el caso de variables aleatorias continuas la definición de un valor esperado es esencialmente la misma, pero las sumatorias se reemplazan con integrales.

Definición 4.1: Sea X una variable aleatoria con distribución de probabilidad $f(x)$. La **media** o **valor esperado** de X es

$$\mu = E(X) = \sum_x xf(x)$$

si X es discreta, y

$$\mu = E(X) = \int_{-\infty}^{\infty} xf(x)dx$$

si X es continua.

El lector debe advertir que la forma para calcular el valor esperado, o media, que se muestra aquí es diferente del método para calcular la media muestral que se describió en el capítulo 1, donde la media muestral se obtuvo usando los datos. En la esperanza matemática el valor esperado se obtiene usando la distribución de probabilidad.

Sin embargo, la media suele considerarse un valor "central" de la distribución subyacente si se utiliza el valor esperado, como en la definición 4.1.

Ejemplo 4.1: Un inspector de calidad obtiene una muestra de un lote que contiene 7 componentes; el lote contiene 4 componentes buenos y 3 defectuosos. El inspector toma una muestra de 3 componentes. Calcule el valor esperado del número de componentes buenos en esta muestra.

Solución: Sea X el número de componentes buenos en la muestra. La distribución de probabilidad de X es

$$f(x) = \frac{\binom{4}{x} \binom{3}{3-x}}{\binom{7}{3}}, \quad x = 0, 1, 2, 3.$$

Unos cálculos sencillos dan $f(0) = 1/35$, $f(1) = 12/35$, $f(2) = 18/35$ y $f(3) = 4/35$. Por lo tanto,

$$\mu = E(X) = (0) \left(\frac{1}{35}\right) + (1) \left(\frac{12}{35}\right) + (2) \left(\frac{18}{35}\right) + (3) \left(\frac{4}{35}\right) = \frac{12}{7} = 1.7.$$

De esta manera, si de un lote de 4 componentes buenos y 3 defectuosos, se seleccionara al azar, una y otra vez, una muestra de tamaño 3, ésta contendría en promedio 1.7 componentes buenos. ■

Ejemplo 4.2: Cierta día un vendedor de una empresa de aparatos médicos tiene dos citas. Considera que en la primera cita tiene 70 por ciento de probabilidades de cerrar una venta, por la cual podría obtener una comisión de \$1000. Por otro lado, cree que en la segunda cita sólo tiene 40 por ciento de probabilidades de cerrar el trato, del cual obtendría \$1500 de comisión. ¿Cuál es su comisión esperada con base en dichas probabilidades? Suponga que los resultados de las citas son independientes.

Solución: En primer lugar sabemos que el vendedor, en las dos citas, puede obtener 4 comisiones totales: \$0, \$1000, \$1500 y \$2500. Necesitamos calcular sus probabilidades asociadas. Mediante la independencia obtenemos

$$\begin{aligned} f(\$0) &= (1 - 0.7)(1 - 0.4) = 0.18, & f(\$2500) &= (0.7)(0.4) = 0.28, \\ f(\$1000) &= (0.7)(1 - 0.4) = 0.42, & \text{y } f(\$1500) &= (1 - 0.7)(0.4) = 0.12. \end{aligned}$$

Por lo tanto, la comisión esperada para el vendedor es

$$\begin{aligned} E(X) &= (\$0)(0.18) + (\$1000)(0.42) + (\$1500)(0.12) + (\$2500)(0.28) \\ &= \$1300. \end{aligned}$$

Los ejemplos 4.1 y 4.2 se diseñaron para que el lector comprenda mejor lo que queremos decir con la frase valor esperado de una variable aleatoria. En ambos casos las variables aleatorias son discretas. Seguimos con un ejemplo de variable aleatoria continua, donde un ingeniero se interesa en la *vida media* de cierto tipo de dispositivo electrónico. Ésta es una ilustración del problema *tiempo que transcurre antes de que se presente una falla* que se enfrenta a menudo en la práctica. El valor esperado de la vida del dispositivo es un parámetro importante para su evaluación. ■

Ejemplo 4.3: Sea X la variable aleatoria que denota la vida en horas de cierto dispositivo electrónico. La función de densidad de probabilidad es

$$f(x) = \begin{cases} \frac{20,000}{x^3}, & x > 100, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la vida esperada para esta clase de dispositivo.

Solución: Si usamos la definición 4.1, tenemos

$$\mu = E(X) = \int_{100}^{\infty} x \frac{20,000}{x^3} dx = \int_{100}^{\infty} \frac{20,000}{x^2} dx = 200.$$

Por lo tanto, esperamos que este tipo de dispositivo dure *en promedio* 200 horas. ▮

Consideremos ahora una nueva variable aleatoria $g(X)$, la cual depende de X ; es decir, cada valor de $g(X)$ es determinado por el valor de X . Por ejemplo, $g(X)$ podría ser X^2 o $3X - 1$, y siempre que X asuma el valor 2, $g(X)$ toma el valor $g(2)$. En particular, si X es una variable aleatoria discreta con distribución de probabilidad $f(x)$, para $x = -1, 0, 1, 2$ y $g(X) = X^2$, entonces,

$$P[g(X) = 0] = P(X = 0) = f(0),$$

$$P[g(X) = 1] = P(X = -1) + P(X = 1) = f(-1) + f(1),$$

$$P[g(X) = 4] = P(X = 2) = f(2),$$

así que la distribución de probabilidad de $g(X)$ se escribe como

$g(x)$	0	1	4
$P[g(X) = g(x)]$	$f(0)$	$f(-1) + f(1)$	$f(2)$

Por medio de la definición del valor esperado de una variable aleatoria obtenemos

$$\begin{aligned} \mu_{g(X)} &= E[g(x)] = 0f(0) + 1[f(-1) + f(1)] + 4f(2) \\ &= (-1)^2f(-1) + (0)^2f(0) + (1)^2f(1) + (2)^2f(2) = \sum_x g(x)f(x). \end{aligned}$$

Este resultado se generaliza en el teorema 4.1 para variables aleatorias discretas y continuas.

Teorema 4.1: Sea X una variable aleatoria con distribución de probabilidad $f(x)$. El valor esperado de la variable aleatoria $g(X)$ es

$$\mu_{g(X)} = E[g(X)] = \sum_x g(x)f(x)$$

si X es discreta, y

$$\mu_{g(X)} = E[g(X)] = \int_{-\infty}^{\infty} g(x)f(x) dx$$

si X es continua.

Ejemplo 4.4: Suponga que el número de automóviles X que pasa por un local de lavado de autos entre las 4:00 P.M. y las 5:00 P.M. de cualquier viernes soleado tiene la siguiente distribución de probabilidad:

x	4	5	6	7	8	9
$P(X = x)$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{6}$	$\frac{1}{6}$

Sea $g(X) = 2X - 1$ la cantidad de dinero en dólares que el administrador paga al operador. Calcule las ganancias esperadas del operador en este periodo específico.

Solución: Por el teorema 4.1, el operador puede esperar recibir

$$\begin{aligned}
 E[g(X)] &= E(2X - 1) = \sum_{x=4}^9 (2x - 1)f(x) \\
 &= (7) \left(\frac{1}{12}\right) + (9) \left(\frac{1}{12}\right) + (11) \left(\frac{1}{4}\right) + (13) \left(\frac{1}{4}\right) \\
 &\quad + (15) \left(\frac{1}{6}\right) + (17) \left(\frac{1}{6}\right) = \$12.67.
 \end{aligned}$$

Ejemplo 4.5: Sea X una variable aleatoria con función de densidad

$$f(x) = \begin{cases} \frac{x^2}{3}, & -1 < x < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el valor esperado de $g(X) = 4X + 3$.

Solución: Por el teorema 4.1 tenemos

$$E(4X + 3) = \int_{-1}^2 \frac{(4x + 3)x^2}{3} dx = \frac{1}{3} \int_{-1}^2 (4x^3 + 3x^2) dx = 8.$$

Debemos extender ahora nuestro concepto de esperanza matemática al caso de dos variables aleatorias X y Y con distribución de probabilidad conjunta $f(x, y)$.

Definición 4.2: Sean X y Y variables aleatorias con distribución de probabilidad conjunta $f(x, y)$. La media o valor esperado de la variable aleatoria $g(X, Y)$ es

$$\mu_{g(X, Y)} = E[g(X, Y)] = \sum_x \sum_y g(x, y)f(x, y)$$

si X y Y son discretas, y

$$\mu_{g(X, Y)} = E[g(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y)f(x, y) dx dy$$

si X y Y son continuas.

Es evidente la generalización de la definición 4.2 para el cálculo de la esperanza matemática de funciones de varias variables aleatorias.

Ejemplo 4.6: Sean X y Y variables aleatorias con la distribución de probabilidad conjunta que se indica en la tabla 3.1 de la página 96. Calcule el valor esperado de $g(X, Y) = XY$. Por conveniencia se repite aquí la tabla.

$f(x, y)$		x			Totales por renglón
		0	1	2	
y	0	$\frac{3}{28}$	$\frac{9}{28}$	$\frac{3}{28}$	$\frac{15}{28}$
	1	$\frac{3}{14}$	$\frac{3}{14}$	0	$\frac{3}{7}$
	2	$\frac{1}{28}$	0	0	$\frac{1}{28}$
Totales por columna		$\frac{5}{14}$	$\frac{15}{28}$	$\frac{3}{28}$	1

Solución: Por la definición 4.2, escribimos

$$\begin{aligned} E(XY) &= \sum_{x=0}^2 \sum_{y=0}^2 xyf(x, y) \\ &= (0)(0)f(0, 0) + (0)(1)f(0, 1) \\ &\quad + (1)(0)f(1, 0) + (1)(1)f(1, 1) + (2)(0)f(2, 0) \\ &= f(1, 1) = \frac{3}{14}. \end{aligned}$$

Ejemplo 4.7: Calcule $E(Y/X)$ para la siguiente función de densidad

$$f(x, y) = \begin{cases} \frac{x(1+3y^2)}{4}, & 0 < x < 2, \quad 0 < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Solución: Tenemos

$$E\left(\frac{Y}{X}\right) = \int_0^1 \int_0^2 \frac{y(1+3y^2)}{4} dx dy = \int_0^1 \frac{y+3y^3}{2} dy = \frac{5}{8}.$$

Observe que si $g(X, Y) = X$ en la definición 4.2, tenemos

$$E(X) = \begin{cases} \sum_x \sum_y xf(x, y) = \sum_x xg(x) & \text{(caso discreto),} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xf(x, y) dy dx = \int_{-\infty}^{\infty} xg(x) dx & \text{(caso continuo),} \end{cases}$$

donde $g(x)$ es la distribución marginal de X . Por lo tanto, para calcular $E(X)$ en un espacio bidimensional, se puede utilizar tanto la distribución de probabilidad conjunta de X y Y , como la distribución marginal de X . De manera similar, definimos

$$E(Y) = \begin{cases} \sum_y \sum_x yf(x, y) = \sum_y yh(y) & \text{(caso discreto),} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} yf(x, y) dx dy = \int_{-\infty}^{\infty} yh(y) dy & \text{(caso continuo),} \end{cases}$$

donde $h(y)$ es la distribución marginal de la variable aleatoria Y .

Ejercicios

4.1 En el ejercicio 3.13 de la página 92 se presenta la siguiente distribución de probabilidad de X , el número de imperfecciones que hay en cada 10 metros de una tela sintética, en rollos continuos de ancho uniforme

x	0	1	2	3	4
$f(x)$	0.41	0.37	0.16	0.05	0.01

Calcule el número promedio de imperfecciones que hay en cada 10 metros de esta tela.

4.2 La distribución de probabilidad de la variable aleatoria discreta X es

$$f(x) = \binom{3}{x} \left(\frac{1}{4}\right)^x \left(\frac{3}{4}\right)^{3-x}, \quad x = 0, 1, 2, 3.$$

Calcule la media de X .

4.3 Calcule la media de la variable aleatoria T que representa el total de las tres monedas del ejercicio 3.25 de la página 93.

4.4 Una moneda está cargada de manera que la probabilidad de ocurrencia de una cara es tres veces mayor que la de una cruz. Calcule el número esperado de cruces si esta moneda se lanza dos veces.

4.5 En un juego de azar a una mujer se le pagan \$3 si saca una jota o una reina, y \$5 si saca un rey o un as de una baraja ordinaria de 52 cartas. Si saca cualquier otra carta, pierde. ¿Cuánto debería pagar si el juego es justo?

4.6 A un operador de un local de lavado de autos se le paga de acuerdo con el número de automóviles que lava. Suponga que las probabilidades de que entre las 4:00 p.m. y las 5:00 p.m. de cualquier viernes soleado reciba \$7, \$9, \$11, \$13, \$15 o \$17 son: $1/12$, $1/12$, $1/4$, $1/4$, $1/6$ y $1/6$, respectivamente. Calcule las ganancias esperadas del operador para este periodo específico.

4.7 Si una persona invierte en unas acciones en particular, en un año tiene una probabilidad de 0.3 de obtener una ganancia de \$4000 o una probabilidad de 0.7 de tener una pérdida de \$1000. ¿Cuál es la ganancia esperada de esta persona?

4.8 Suponga que un distribuidor de joyería antigua está interesado en comprar un collar de oro para el que tiene 0.22 de probabilidades de venderlo con \$250 de utilidad; 0.36 de venderlo con \$150 de utilidad; 0.28 de venderlo al costo y 0.14 de venderlo con una pérdida de \$150. ¿Cuál es su utilidad esperada?

4.9 Un piloto privado desea asegurar su avión por \$200,000. La aseguradora estima que la probabilidad de pérdida total es de 0.002, que la probabilidad de una pérdida del 50% es de 0.01 y la probabilidad de una

pérdida del 25% es de 0.1. Si se ignoran todas las demás pérdidas parciales, ¿qué prima debería cobrar cada año la aseguradora para tener una utilidad promedio de \$500?

4.10 Dos expertos en calidad de neumáticos examinan lotes de éstos y asignan a cada neumático puntuaciones de calidad en una escala de tres puntos. Sea X la puntuación dada por el experto A y Y la dada por el experto B . La siguiente tabla presenta la distribución conjunta para X y Y .

$f(x, y)$		y		
		1	2	3
x	1	0.10	0.05	0.02
	2	0.10	0.35	0.05
	3	0.03	0.10	0.20

Calcule μ_x y μ_y .

4.11 La función de densidad de las mediciones codificadas del diámetro de paso de los hilos de un encaje es

$$f(x) = \begin{cases} \frac{4}{\pi(1+x^2)}, & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el valor esperado de X .

4.12 Si la utilidad para un distribuidor de un automóvil nuevo, en unidades de \$5000, se puede ver como una variable aleatoria X que tiene la siguiente función de densidad

$$f(x) = \begin{cases} 2(1-x), & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la utilidad promedio por automóvil.

4.13 La función de densidad de la variable aleatoria continua X , el número total de horas que una familia utiliza una aspiradora durante un año, en unidades de 100 horas, se da en el ejercicio 3.7 de la página 92 como

$$f(x) = \begin{cases} x, & 0 < x < 1, \\ 2-x, & 1 \leq x < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el número promedio de horas por año que las familias utilizan sus aspiradoras.

4.14 Calcule la proporción X de personas que se podría esperar que respondieran a cierta encuesta que se envía por correo, si X tiene la siguiente función de densidad

$$f(x) = \begin{cases} \frac{2(x+2)}{5}, & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

4.15 Suponga que dos variables aleatorias (X, Y) están distribuidas de manera uniforme en un círculo con radio a . Entonces, la función de densidad de probabilidad conjunta es

$$f(x, y) = \begin{cases} \frac{1}{\pi a^2}, & x^2 + y^2 \leq a^2, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule μ_x , el valor esperado de X .

4.16 Suponga que usted inspecciona un lote de 1000 bombillas de luz, entre las cuales hay 20 defectuosas, y elige al azar dos bombillas del lote sin reemplazo. Sean

$$X_1 = \begin{cases} 1, & \text{si la primera bombilla está defectuosa,} \\ 0, & \text{en otro caso.} \end{cases}$$

$$X_2 = \begin{cases} 1, & \text{si la segunda bombilla está defectuosa,} \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la probabilidad de que al menos una de las bombillas elegidas esté defectuosa. [Sugerencia: Calcule $P(X_1 + X_2 = 1)$.]

4.17 Sea X una variable aleatoria con la siguiente distribución de probabilidad:

x	-3	6	9
$f(x)$	1/6	1/2	1/3

Calcule $\mu_{g(X)}$, donde $g(X) = (2X + 1)^2$.

4.18 Calcule el valor esperado de la variable aleatoria $g(X) = X^2$, donde X tiene la distribución de probabilidad del ejercicio 4.2.

4.19 Una empresa industrial grande compra varios procesadores de textos nuevos al final de cada año; el número exacto depende de la frecuencia de reparaciones del año anterior. Suponga que el número de procesadores de textos, X , que se compran cada año tiene la siguiente distribución de probabilidad:

x	0	1	2	3
$f(x)$	1/10	3/10	2/5	1/5

Si el costo del modelo deseado es de \$1200 por unidad y al final del año la empresa obtiene un descuento de $50X^2$ dólares, ¿cuánto espera gastar esta empresa en nuevos procesadores de textos durante este año?

4.20 Una variable aleatoria continua X tiene la siguiente función de densidad

$$f(x) = \begin{cases} e^{-x}, & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el valor esperado de $g(X) = e^{2X/3}$.

4.21 ¿Cuál es la utilidad promedio por automóvil que obtiene un distribuidor, si la utilidad en cada uno está dada por $g(X) = X^2$, donde X es una variable aleatoria que tiene la función de densidad del ejercicio 4.12?

4.22 El periodo de hospitalización, en días, para pacientes que siguen el tratamiento para cierto tipo de trastorno renal es una variable aleatoria $Y = X + 4$, donde X tiene la siguiente función de densidad

$$f(x) = \begin{cases} \frac{32}{(x+4)^3}, & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el número promedio de días que una persona permanece hospitalizada con el fin de seguir el tratamiento para dicha enfermedad.

4.23 Suponga que X y Y tienen la siguiente función de probabilidad conjunta:

$f(x, y)$		x	
		2	4
y	1	0.10	0.15
	3	0.20	0.30
	5	0.10	0.15

- a) Calcule el valor esperado de $g(X, Y) = XY^2$.
b) Calcule μ_x y μ_y .

4.24 Remítase a las variables aleatorias cuya distribución de probabilidad conjunta se da en el ejercicio 3.39 de la página 105 y

- a) calcule $E(X^2Y - 2XY)$;
b) calcule $\mu_x - \mu_y$.

4.25 Remítase a las variables aleatorias cuya distribución de probabilidad conjunta se da en el ejercicio 3.51 de la página 106 y calcule la media para el número total de jotas y reyes cuando se sacan 3 cartas, sin reemplazo, de las 12 cartas mayores de una baraja ordinaria de 52 cartas.

4.26 Sean X y Y las siguientes variables aleatorias con función de densidad conjunta

$$f(x, y) = \begin{cases} 4xy, & 0 < x, y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el valor esperado de $Z = \sqrt{X^2 + Y^2}$.

4.27 En el ejercicio 3.27 de la página 93 una función de densidad está dada por el tiempo que tarda en fallar un componente importante de un reproductor de DVD. Calcule el número medio de horas antes de que empiece a fallar el componente y, por lo tanto, el reproductor de DVD.

4.28 Considere la información del ejercicio 3.28 de la página 93. El problema tiene que ver con el peso, en onzas, del producto que contiene una caja de cereal con

$$f(x) = \begin{cases} \frac{2}{3}, & 23.75 \leq x \leq 26.25, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Grafique la función de densidad.
 b) Calcule el valor esperado o peso medio en onzas.
 c) ¿Se sorprende de su respuesta en b)? Explique lo que responda.

4.29 El ejercicio 3.29 de la página 93 se refiere a una importante distribución del tamaño de las partículas caracterizada por

$$f(x) = \begin{cases} 3x^{-4}, & x > 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Grafique la función de densidad.
 b) Determine el tamaño medio de la partícula.

4.30 En el ejercicio 3.31 de la página 94 la distribución del tiempo que transcurre antes de que una lavadora requiera una reparación mayor fue dada como

$$f(y) = \begin{cases} \frac{1}{4}e^{-y/4}, & y \geq 0, \\ 0, & \text{en otro caso.} \end{cases}$$

¿Cuál es la media de población del tiempo que transcurre antes de requerir la reparación?

4.31 Considere el ejercicio 3.32 de la página 94.

- a) ¿Cuál es la proporción media del presupuesto asignado para el control ambiental y de la contaminación?
 b) ¿Cuál es la probabilidad de que una empresa elegida al azar tenga una proporción asignada para el control ambiental y de la contaminación que exceda la media de la población dada en a)?

4.32 En el ejercicio 3.13 de la página 92 la distribución del número de imperfecciones en cada 10 metros de tela sintética fue dada por

x	0	1	2	3	4
$f(x)$	0.41	0.37	0.16	0.05	0.01

- a) Grafique la función de probabilidad.
 b) Calcule el número de imperfecciones esperado $E(X) = \mu$.
 c) Calcule $E(X^2)$.

4.2 Varianza y covarianza de variables aleatorias

La media o valor esperado de una variable aleatoria X es de especial importancia en estadística porque describe en dónde se centra la distribución de probabilidad. Sin embargo, la media por sí misma no ofrece una descripción adecuada de la forma de la distribución. También se necesita clasificar la variabilidad en la distribución. En la figura 4.1 tenemos los histogramas de dos distribuciones de probabilidad discretas con la misma media $\mu = 2$, pero que difieren de manera considerable en la variabilidad o dispersión de sus observaciones sobre la media.

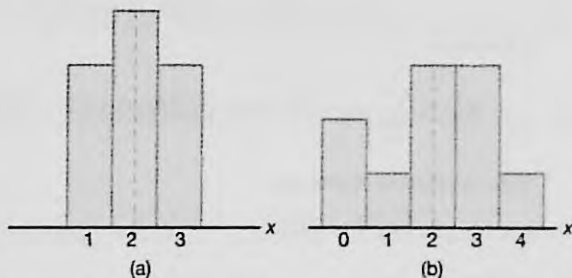


Figura 4.1: Distribuciones con medias iguales y dispersiones diferentes.

La medida de variabilidad más importante de una variable aleatoria X se obtiene aplicando el teorema 4.1 con $g(X) = (X - \mu)^2$. A esta cantidad se le denomina **varianza** de la variable aleatoria X o **varianza** de la distribución de probabilidad de X y se

denota como $\text{Var}(X)$, o con el símbolo σ_X^2 , o simplemente como σ^2 cuando es evidente a qué variable aleatoria se está haciendo referencia.

Definición 4.3: Sea X una variable aleatoria con distribución de probabilidad $f(x)$ y media μ . La varianza de X es

$$\sigma^2 = E[(X - \mu)^2] = \sum_x (x - \mu)^2 f(x), \quad \text{si } X \text{ es discreta, y}$$

$$\sigma^2 = E[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx, \quad \text{si } X \text{ es continua.}$$

La raíz cuadrada positiva de la varianza, σ , se llama **desviación estándar** de X .

La cantidad $x - \mu$ en la definición 4.3 se llama **desviación de una observación** respecto a su media. Como estas desviaciones se elevan al cuadrado y después se promedian, σ^2 será mucho menor para un conjunto de valores x que estén cercanos a μ , que para un conjunto de valores que varíe de forma considerable de μ .

Ejemplo 4.8: Suponga que la variable aleatoria X representa el número de automóviles que se utilizan con propósitos de negocios oficiales en un día de trabajo dado. La distribución de probabilidad para la empresa A [figura 4.1(a)] es

x	1	2	3
$f(x)$	0.3	0.4	0.3

y para la empresa B [figura 4.1(b)] es

x	0	1	2	3	4
$f(x)$	0.2	0.1	0.3	0.3	0.1

Demuestre que la varianza de la distribución de probabilidad para la empresa B es mayor que la de la empresa A .

Solución: Para la empresa A encontramos que

$$\mu_A = E(X) = (1)(0.3) + (2)(0.4) + (3)(0.3) = 2.0,$$

y entonces

$$\sigma_A^2 = \sum_{x=1}^3 (x - 2)^2 = (1 - 2)^2(0.3) + (2 - 2)^2(0.4) + (3 - 2)^2(0.3) = 0.6.$$

Para la empresa B tenemos

$$\mu_B = E(X) = (0)(0.2) + (1)(0.1) + (2)(0.3) + (3)(0.3) + (4)(0.1) = 2.0,$$

y entonces

$$\begin{aligned} \sigma_B^2 &= \sum_{x=0}^4 (x - 2)^2 f(x) \\ &= (0 - 2)^2(0.2) + (1 - 2)^2(0.1) + (2 - 2)^2(0.3) \\ &\quad + (3 - 2)^2(0.3) + (4 - 2)^2(0.1) = 1.6. \end{aligned}$$

Es evidente que la varianza del número de automóviles que se utilizan con propósitos de negocios oficiales es mayor para la empresa *B* que para la empresa *A*. ■

Una fórmula alternativa que se prefiere para calcular σ^2 , que a menudo simplifica los cálculos, se establece en el siguiente teorema.

Teorema 4.2: La varianza de una variable aleatoria X es

$$\sigma^2 = E(X^2) - \mu^2.$$

Prueba: Para el caso discreto escribimos

$$\begin{aligned}\sigma^2 &= \sum_x (x - \mu)^2 f(x) = \sum_x (x^2 - 2\mu x + \mu^2) f(x) \\ &= \sum_x x^2 f(x) - 2\mu \sum_x x f(x) + \mu^2 \sum_x f(x).\end{aligned}$$

Como $\mu = \sum_x x f(x)$ por definición, y $\sum_x f(x) = 1$ para cualquier distribución de probabilidad discreta, se deduce que

$$\sigma^2 = \sum_x x^2 f(x) - \mu^2 = E(X^2) - \mu^2.$$

Para el caso continuo la demostración es la misma paso a paso, reemplazando las sumatorias por integrales. ■

Ejemplo 4.9: Suponga que la variable aleatoria X representa el número de partes defectuosas de una máquina cuando de una línea de producción se obtiene una muestra de tres partes y se somete a prueba. La siguiente es la distribución de probabilidad de X .

x	0	1	2	3
$f(x)$	0.51	0.38	0.10	0.01

Utilice el teorema 4.2 y calcule σ^2 .

Solución: Primero calculamos

$$\mu = (0)(0.51) + (1)(0.38) + (2)(0.10) + (3)(0.01) = 0.61.$$

Luego,

$$E(X^2) = (0)(0.51) + (1)(0.38) + (4)(0.10) + (9)(0.01) = 0.87.$$

Por lo tanto,

$$\sigma^2 = 0.87 - (0.61)^2 = 0.4979. \quad \blacksquare$$

Ejemplo 4.10: La demanda semanal de una bebida para una cadena local de tiendas de abarrotes, en miles de litros, es una variable aleatoria continua X que tiene la siguiente densidad de probabilidad

$$f(x) = \begin{cases} 2(x-1), & 1 < x < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la media y la varianza de X .

Solución: Al calcular $E(X)$ y $E(X^2)$ tenemos

$$\mu = E(X) = 2 \int_1^2 x(x-1) dx = \frac{5}{3}$$

y

$$E(X^2) = 2 \int_1^2 x^2(x-1) dx = \frac{17}{6}.$$

Por lo tanto,

$$\sigma^2 = \frac{17}{6} - \left(\frac{5}{3}\right)^2 = \frac{1}{18}.$$

Hasta el momento la varianza o la desviación estándar sólo tiene significado cuando comparamos dos o más distribuciones que tienen las mismas unidades de medida. Por lo tanto, podemos comparar las varianzas de las distribuciones de contenido, medido en litros, de botellas de jugo de naranja de dos empresas, y el valor más grande indicaría la empresa cuyo producto es más variable o menos uniforme. No tendría caso comparar la varianza de una distribución de estaturas con la varianza de una distribución de calificaciones de aptitud. En la sección 4.4 mostramos cómo se utiliza la desviación estándar para describir una sola distribución de observaciones.

Extenderemos ahora nuestro concepto de varianza de una variable aleatoria X para incluir también variables aleatorias relacionadas con X . Para la variable aleatoria $g(X)$ la varianza se denotará por $\sigma_{g(X)}^2$ y se calculará empleando el siguiente teorema.

Teorema 4.3: Sea X una variable aleatoria con distribución de probabilidad $f(x)$. La varianza de la variable aleatoria $g(X)$ es

$$\sigma_{g(X)}^2 = E \{ [g(X) - \mu_{g(X)}]^2 \} = \sum_x [g(x) - \mu_{g(X)}]^2 f(x)$$

si X es discreta, y

$$\sigma_{g(X)}^2 = E \{ [g(X) - \mu_{g(X)}]^2 \} = \int_{-\infty}^{\infty} [g(x) - \mu_{g(X)}]^2 f(x) dx$$

si X es continua.

Prueba: Como $g(X)$ es en sí misma una variable aleatoria con media $\mu_{g(X)}$, como se define en el teorema 4.1, de la definición 4.3 se deduce que

$$\sigma_{g(X)}^2 = E \{ [g(X) - \mu_{g(X)}]^2 \}.$$

Ahora bien, la demostración se completa aplicando nuevamente el teorema 4.1 a la variable aleatoria $[g(X) - \mu_{g(X)}]^2$.

Ejemplo 4.11: Calcule la varianza de $g(X) = 2X + 3$, donde X es una variable aleatoria con la siguiente distribución de probabilidad

x	0	1	2	3
$f(x)$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{2}$	$\frac{1}{8}$

Solución: Primero se calcula la media de la variable aleatoria $2X + 3$. De acuerdo con el teorema 4.1,

$$\mu_{2X+3} = E(2X + 3) = \sum_{x=0}^3 (2x + 3)f(x) = 6.$$

Ahora, usando el teorema 4.3, tenemos

$$\begin{aligned}\sigma_{2X+3}^2 &= E\{[(2X + 3) - \mu_{2X+3}]^2\} = E\{(2X + 3 - 6)^2\} \\ &= E(4X^2 - 12X + 9) = \sum_{x=0}^3 (4x^2 - 12x + 9)f(x) = 4.\end{aligned}$$

Ejemplo 4.12: Sea X una variable aleatoria que tiene la función de densidad dada en el ejemplo 4.5 de la página 115. Calcule la varianza de la variable aleatoria $g(X) = 4X + 3$.

Solución: En el ejemplo 4.5 encontramos que $\mu_{4X+3} = 8$. Ahora bien, usando el teorema 4.3,

$$\begin{aligned}\sigma_{4X+3}^2 &= E\{[(4X + 3) - 8]^2\} = E\{(4X - 5)^2\} \\ &= \int_{-1}^2 (4x - 5)^2 \frac{x^2}{3} dx = \frac{1}{3} \int_{-1}^2 (16x^4 - 40x^3 + 25x^2) dx = \frac{51}{5}.\end{aligned}$$

Si $g(X, Y) = (X - \mu_X)(Y - \mu_Y)$, donde $\mu_X = E(X)$ y $\mu_Y = E(Y)$, la definición 4.2 da un valor esperado denominado **covarianza** de X y Y , que se denota por σ_{XY} o $\text{Cov}(X, Y)$.

Definición 4.4: Sean X y Y variables aleatorias con distribución de probabilidad conjunta $f(x, y)$. La covarianza de X y Y es

$$\sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)] = \sum_x \sum_y (x - \mu_X)(y - \mu_Y)f(x, y)$$

si X y Y son discretas, y

$$\sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_X)(y - \mu_Y)f(x, y) dx dy$$

si X y Y son continuas.

La covarianza entre dos variables aleatorias es una medida de la naturaleza de la asociación entre ambas. Si valores grandes de X a menudo dan como resultado valores grandes de Y , o valores pequeños de X , dan como resultado valores pequeños de Y , $X - \mu_X$ positiva con frecuencia dará como resultado $Y - \mu_Y$ positiva, y $X - \mu_X$ negativa a menudo dará como resultado $Y - \mu_Y$ negativa. Por consiguiente, el producto $(X - \mu_X)(Y - \mu_Y)$ tenderá a ser positivo. Por otro lado, si con frecuencia valores grandes de X dan como resultado valores pequeños de Y , entonces el producto $(X - \mu_X)(Y - \mu_Y)$ tenderá a ser negativo. El *signo* de la covarianza indica si la relación entre dos variables aleatorias dependientes es positiva o negativa. Cuando X y Y son estadísticamente independientes, se puede demostrar que la covarianza es cero (véase el corolario 4.5). Lo opuesto, sin embargo, por lo general no es cierto. Dos variables pueden tener covarianza cero y aun así no ser estadísticamente independientes. Observe que la covarianza sólo describe la relación *lineal* entre dos variables aleatorias. Por consiguiente, si una covarianza entre X y Y es cero, X y Y podrían tener una relación no lineal, lo cual significa que no necesariamente son independientes.

La fórmula alternativa que se prefiere para σ_{XY} se establece en el teorema 4.4.

Teorema 4.4: La covarianza de dos variables aleatorias X y Y , con medias μ_x y μ_y , respectivamente, está dada por

$$\sigma_{XY} = E(XY) - \mu_x \mu_y.$$

Prueba: Para el caso discreto escribimos

$$\begin{aligned}\sigma_{XY} &= \sum_x \sum_y (x - \mu_x)(y - \mu_y) f(x, y) \\ &= \sum_x \sum_y xy f(x, y) - \mu_x \sum_x \sum_y y f(x, y) \\ &\quad - \mu_y \sum_x \sum_y x f(x, y) + \mu_x \mu_y \sum_x \sum_y f(x, y).\end{aligned}$$

Dado que

$$\mu_x = \sum_x x f(x, y), \quad \mu_y = \sum_y y f(x, y), \quad \text{y} \quad \sum_x \sum_y f(x, y) = 1$$

para cualquier distribución discreta conjunta se deduce que

$$\sigma_{XY} = E(XY) - \mu_x \mu_y - \mu_y \mu_x + \mu_x \mu_y = E(XY) - \mu_x \mu_y.$$

Para el caso continuo la demostración es idéntica, pero las sumatorias se reemplazan por integrales.

Ejemplo 4.13: En el ejemplo 3.14 de la página 95 se describe una situación acerca del número de repuestos azules X y el número de repuestos rojos Y . Cuando de cierta caja se seleccionan dos repuestos para bolígrafo al azar y la distribución de probabilidad conjunta es la siguiente,

$f(x, y)$		x			$h(y)$
		0	1	2	
y	0	$\frac{3}{28}$	$\frac{9}{28}$	$\frac{3}{28}$	$\frac{15}{28}$
	1	$\frac{3}{14}$	$\frac{3}{14}$	0	$\frac{3}{7}$
	2	$\frac{1}{28}$	0	0	$\frac{1}{28}$
$g(x)$		$\frac{5}{14}$	$\frac{15}{28}$	$\frac{3}{28}$	1

Calcule la covarianza de X y Y .

Solución: Del ejemplo 4.6 vemos que $E(XY) = 3/14$. Ahora bien,

$$\mu_x = \sum_{x=0}^2 xg(x) = (0) \left(\frac{5}{14}\right) + (1) \left(\frac{15}{28}\right) + (2) \left(\frac{3}{28}\right) = \frac{3}{4},$$

y

$$\mu_y = \sum_{y=0}^2 yh(y) = (0) \left(\frac{15}{28}\right) + (1) \left(\frac{3}{7}\right) + (2) \left(\frac{1}{28}\right) = \frac{1}{2}.$$

Por lo tanto,

$$\sigma_{xy} = E(XY) - \mu_x \mu_y = \frac{3}{14} - \left(\frac{3}{4}\right) \left(\frac{1}{2}\right) = -\frac{9}{56}.$$

Ejemplo 4.14: La fracción X de corredores y la fracción Y de corredoras que compiten en carreras de maratón se describen mediante la función de densidad conjunta

$$f(x, y) = \begin{cases} 8xy, & 0 \leq y \leq x \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la covarianza de X y Y .

Solución: Primero calculamos las funciones de densidad marginal. Éstas son

$$g(x) = \begin{cases} 4x^3, & 0 \leq x \leq 1, \\ 0, & \text{en otro caso,} \end{cases}$$

y

$$h(y) = \begin{cases} 4y(1-y^2), & 0 \leq y \leq 1, \\ 0, & \text{en otro caso,} \end{cases}$$

A partir de las funciones de densidad marginal dadas, calculamos

$$\mu_x = E(X) = \int_0^1 4x^4 dx = \frac{4}{5} \quad \text{y} \quad \mu_y = \int_0^1 4y^2(1-y^2) dy = \frac{8}{15}.$$

De las funciones de densidad conjunta dadas arriba, tenemos

$$E(XY) = \int_0^1 \int_y^1 8x^2y^2 dx dy = \frac{4}{9}.$$

Entonces,

$$\sigma_{xy} = E(XY) - \mu_x \mu_y = \frac{4}{9} - \left(\frac{4}{5}\right) \left(\frac{8}{15}\right) = \frac{4}{225}. \quad \blacksquare$$

Aunque la covarianza entre dos variables aleatorias brinda información respecto de la naturaleza de la relación, la magnitud de σ_{xy} no indica nada respecto a la fuerza de la relación, ya que σ_{xy} depende de la escala. Su magnitud dependerá de las unidades que se utilicen para medir X y Y . Hay una versión de la covarianza sin escala que se denomina **coeficiente de correlación** y se utiliza ampliamente en estadística.

Definición 4.5: Sean X y Y variables aleatorias con covarianza σ_{xy} y desviaciones estándar σ_x y σ_y , respectivamente. El coeficiente de correlación de X y Y es

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}.$$

Debería quedar claro para el lector que ρ_{xy} no tiene las unidades de X y Y . El coeficiente de correlación satisface la desigualdad $-1 \leq \rho_{xy} \leq 1$. Toma un valor de cero cuando $\sigma_{xy} = 0$. Donde hay una dependencia lineal exacta, digamos $Y \equiv a + bX$, $\rho_{xy} = 1$ si

$b > 0$ y $\rho_{XY} = -1$ si $b < 0$. (Véase el ejercicio 4.48). En el capítulo 12, donde examinaremos la regresión lineal, analizamos más a fondo el coeficiente de correlación.

Ejemplo 4.15: Calcule el coeficiente de correlación entre X y Y en el ejemplo 4.13.

Solución: Dado que

$$E(X^2) = (0^2) \left(\frac{5}{14}\right) + (1^2) \left(\frac{15}{28}\right) + (2^2) \left(\frac{3}{28}\right) = \frac{27}{28}$$

y

$$E(Y^2) = (0^2) \left(\frac{15}{28}\right) + (1^2) \left(\frac{3}{7}\right) + (2^2) \left(\frac{1}{28}\right) = \frac{4}{7},$$

obtenemos

$$\sigma_X^2 = \frac{27}{28} - \left(\frac{3}{4}\right)^2 = \frac{45}{112} \quad \text{y} \quad \sigma_Y^2 = \frac{4}{7} - \left(\frac{1}{2}\right)^2 = \frac{9}{28}.$$

Por lo tanto, el coeficiente de correlación entre X y Y es

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} = \frac{-9/56}{\sqrt{(45/112)(9/28)}} = -\frac{1}{\sqrt{5}}.$$

Ejemplo 4.16: Calcule el coeficiente de correlación entre X y Y en el ejemplo 4.14.

Solución: Dado que

$$E(X^2) = \int_0^1 4x^5 dx = \frac{2}{3} \quad \text{y} \quad E(Y^2) = \int_0^1 4y^3(1-y^2) dy = 1 - \frac{2}{3} = \frac{1}{3},$$

concluimos que

$$\sigma_X^2 = \frac{2}{3} - \left(\frac{4}{5}\right)^2 = \frac{2}{75} \quad \text{y} \quad \sigma_Y^2 = \frac{1}{3} - \left(\frac{8}{15}\right)^2 = \frac{11}{225}.$$

Por lo tanto,

$$\rho_{XY} = \frac{4/225}{\sqrt{(2/75)(11/225)}} = \frac{4}{\sqrt{66}}.$$

Observe que, aunque la covarianza en el ejemplo 4.15 tiene mayor magnitud (sin importar el signo) que la del ejemplo 4.16, la relación entre las magnitudes de los coeficientes de correlación en estos dos ejemplos es exactamente la inversa. Esto es evidencia de que no debemos basarnos en la magnitud de la covarianza para determinar la fuerza de la relación.

Ejercicios

4.33 Use la definición 4.3 de la página 120 para encontrar la varianza de la variable aleatoria X del ejercicio 4.7 de la página 117.

4.34 Sea X una variable aleatoria con la siguiente distribución de probabilidad:

x	-2	3	5
$f(x)$	0.3	0.2	0.5

Calcule la desviación estándar de X .

4.35 La variable aleatoria X , que representa el número de errores por 100 líneas de código de programación, tiene la siguiente distribución de probabilidad:

x	2	3	4	5	6
$f(x)$	0.01	0.25	0.4	0.3	0.04

Utilice el teorema 4.2 de la página 121 para calcular la varianza de X .

4.36 Suponga que las probabilidades de que 0, 1, 2 o 3 fallas de energía eléctrica afecten cierta subdivisión en cualquier año dado son 0.4, 0.3, 0.2 y 0.1, respectivamente. Calcule la media y la varianza de la variable aleatoria X que representa el número de fallas de energía que afectan esta subdivisión.

4.37 La utilidad que obtiene un distribuidor, en unidades de \$5000, al vender un automóvil nuevo es una variable aleatoria X que tiene la función de densidad que se presenta en el ejercicio 4.12 de la página 117. Calcule la varianza de X .

4.38 La proporción de personas que responden cierta encuesta que se manda por correo es una variable aleatoria X , la cual tiene la función de densidad del ejercicio 4.14 de la página 117. Calcule la varianza de X .

4.39 El número total de horas que una familia utiliza una aspiradora en un año, en unidades de 100 horas, es una variable aleatoria X que tiene la función de densidad dada en el ejercicio 4.13 de la página 117. Calcule la varianza de X .

4.40 Remítase al ejercicio 4.14 de la página 117 y calcule $\sigma_{g(X)}^2$ para la función $g(X) = 3X^2 + 4$.

4.41 Calcule la desviación estándar de la variable aleatoria $g(X) = (2X + 1)^2$ del ejercicio 4.17 en la página 118.

4.42 Utilice los resultados del ejercicio 4.21 de la página 118 y calcule la varianza de $g(X) = X^2$, donde X es una variable aleatoria que tiene la función de densidad del ejercicio 4.12 de la página 117.

4.43 El tiempo que transcurre, en minutos, para que un avión obtenga vía libre para despegar en cierto aeropuerto es una variable aleatoria $Y = 3X - 2$, donde X tiene la siguiente función de densidad

$$f(x) = \begin{cases} \frac{1}{4}e^{-x/4}, & x > 0 \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la media y la varianza de la variable aleatoria Y .

4.44 Calcule la covarianza de las variables aleatorias X y Y del ejercicio 3.39 de la página 105.

4.45 Calcule la covarianza de las variables aleatorias X y Y del ejercicio 3.49 de la página 106.

4.46 Calcule la covarianza de las variables aleatorias X y Y del ejercicio 3.44 de la página 105.

4.47 Calcule la covarianza de las variables aleatorias X y Y cuya función de densidad conjunta está dada en el ejercicio 3.40 de la página 105.

4.48 Dada una variable aleatoria X , con desviación estándar σ_x y una variable aleatoria $Y = a + bX$, demuestre que si $b < 0$, el coeficiente de correlación $\rho_{XY} = -1$, y si $b > 0$, $\rho_{XY} = 1$.

4.49 Considere la situación del ejercicio 4.32 de la página 119. La distribución del número de imperfecciones por cada 10 metros de tela sintética está dada por

x	0	1	2	3	4
$f(x)$	0.41	0.37	0.16	0.05	0.01

Calcule la varianza y la desviación estándar del número de imperfecciones.

4.50 En una tarea de laboratorio, si el equipo está funcionando, la función de densidad del resultado observado X es

$$f(x) = \begin{cases} 2(1-x), & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule la varianza y la desviación estándar de X .

4.51 Determine el coeficiente de correlación entre X y Y para las variables aleatorias X y Y del ejercicio 3.39 de la página 105.

4.52 Las variables aleatorias X y Y tienen la siguiente distribución conjunta

$$f(x, y) = \begin{cases} 2, & 0 < x \leq y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Determine el coeficiente de correlación entre X y Y .

4.3 Medias y varianzas de combinaciones lineales de variables aleatorias

Ahora estudiaremos algunas propiedades útiles que simplificarán los cálculos de las medias y las varianzas de variables aleatorias que aparecen en los siguientes capítulos. Estas propiedades nos permitirán ocuparnos de las esperanzas matemáticas en términos de otros parámetros que ya conocemos o que ya calculamos con facilidad. Todos los resultados que presentamos aquí son válidos para variables aleatorias continuas y discretas. Las demostraciones se dan sólo para el caso continuo. Comenzamos con un teorema y dos corolarios que deberían ser, de forma intuitiva, razonables para el lector.

Teorema 4.5: Si a y b son constantes, entonces,

$$E(aX + b) = aE(X) + b.$$

Prueba: Por la definición de valor esperado,

$$E(aX + b) = \int_{-\infty}^{\infty} (ax + b)f(x) dx = a \int_{-\infty}^{\infty} xf(x) dx + b \int_{-\infty}^{\infty} f(x) dx.$$

La primera integral de la derecha es $E(X)$ y la segunda integral es igual a 1. Por lo tanto,

$$E(aX + b) = aE(X) + b. \quad \square$$

Corolario 4.1: Al establecer que $a = 0$ vemos que $E(b) = b$.

Corolario 4.2: Al establecer que $b = 0$ vemos que $E(aX) = aE(X)$.

Ejemplo 4.17: Aplique el teorema 4.5 a la variable aleatoria discreta $f(X) = 2X - 1$ para resolver de nuevo el ejemplo 4.4 de la página 115.

Solución: De acuerdo con el teorema 4.5, escribimos

$$E(2X - 1) = 2E(X) - 1.$$

Ahora,

$$\begin{aligned} \mu &= E(X) = \sum_{x=4}^9 xf(x) \\ &= (4) \left(\frac{1}{12}\right) + (5) \left(\frac{1}{12}\right) + (6) \left(\frac{1}{4}\right) + (7) \left(\frac{1}{4}\right) + (8) \left(\frac{1}{6}\right) + (9) \left(\frac{1}{6}\right) = \frac{41}{6}. \end{aligned}$$

Por lo tanto,

$$\mu_{2X-1} = (2) \left(\frac{41}{6}\right) - 1 = \$12.67,$$

como antes.

Ejemplo 4.18: Para resolver de nuevo el ejemplo 4.5 de la página 115 aplique el teorema 4.5 a la variable aleatoria continua $g(X) = 4X + 3$.

Solución: En el ejemplo 4.5 utilizamos el teorema 4.5 para escribir

$$E(4X + 3) = 4E(X) + 3.$$

Ahora,

$$E(X) = \int_{-1}^2 x \left(\frac{x^2}{3}\right) dx = \int_{-1}^2 \frac{x^3}{3} dx = \frac{5}{4}.$$

Por lo tanto,

$$E(4X + 3) = (4) \left(\frac{5}{4}\right) + 3 = 8,$$

como antes. ▮

Teorema 4.6: El valor esperado de la suma o diferencia de dos o más funciones de una variable aleatoria X es la suma o diferencia de los valores esperados de las funciones. Es decir,

$$E[g(X) \pm h(X)] = E[g(X)] \pm E[h(X)].$$

Prueba: Por definición,

$$\begin{aligned} E[g(X) \pm h(X)] &= \int_{-\infty}^{\infty} [g(x) \pm h(x)]f(x) dx \\ &= \int_{-\infty}^{\infty} g(x)f(x) dx \pm \int_{-\infty}^{\infty} h(x)f(x) dx \\ &= E[g(X)] \pm E[h(X)]. \end{aligned}$$
▮

Ejemplo 4.19: Sea X una variable aleatoria con la siguiente distribución de probabilidad:

x	0	1	2	3
$f(x)$	$\frac{1}{3}$	$\frac{1}{2}$	0	$\frac{1}{6}$

Calcule el valor esperado de $Y = (X - 1)^2$.

Solución: Si aplicamos el teorema 4.6 a la función $Y = (X - 1)^2$, podemos escribir

$$E[(X - 1)^2] = E(X^2 - 2X + 1) = E(X^2) - 2E(X) + E(1).$$

A partir del corolario 4.1, $E(1) = 1$, y por cálculo directo

$$E(X) = (0) \left(\frac{1}{3}\right) + (1) \left(\frac{1}{2}\right) + (2)(0) + (3) \left(\frac{1}{6}\right) = 1 \text{ y}$$

$$E(X^2) = (0) \left(\frac{1}{3}\right) + (1) \left(\frac{1}{2}\right) + (4)(0) + (9) \left(\frac{1}{6}\right) = 2.$$

En consecuencia,

$$E[(X - 1)^2] = 2 - (2)(1) + 1 = 1. \quad \text{▮}$$

Ejemplo 4.20: La demanda semanal de cierta bebida en una cadena de tiendas de abarrotes, en miles de litros, es una variable aleatoria continua $g(X) = X^2 + X - 2$, donde X tiene la siguiente función de densidad

$$f(x) = \begin{cases} 2(x-1), & 1 < x < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el valor esperado para la demanda semanal de la bebida.

Solución: Por medio del teorema 4.6, escribimos

$$E(X^2 + X - 2) = E(X^2) + E(X) - E(2).$$

A partir del corolario 4.1, $E(2) = 2$, y por integración directa,

$$E(X) = \int_1^2 2x(x-1) dx = \frac{5}{3} \quad \text{y} \quad E(X^2) = \int_1^2 2x^2(x-1) dx = \frac{17}{6}.$$

Entonces,

$$E(X^2 + X - 2) = \frac{17}{6} + \frac{5}{3} - 2 = \frac{5}{2},$$

así que la demanda semanal promedio de la bebida en esta cadena de tiendas de abarrotes es de 2500 litros. J

Suponga que tenemos dos variables aleatorias X y Y con distribución de probabilidad conjunta $f(x, y)$. Dos propiedades adicionales que serán muy útiles en los capítulos siguientes incluyen los valores esperados de la suma, la diferencia y el producto de estas dos variables aleatorias. Sin embargo, comenzaremos por demostrar un teorema sobre el valor esperado de la suma o diferencia de funciones de las variables dadas. Por supuesto, tan sólo se trata de una extensión del teorema 4.6.

Teorema 4.7: El valor esperado de la suma o diferencia de dos o más funciones de las variables aleatorias X y Y es la suma o diferencia de los valores esperados de las funciones. Es decir,

$$E[g(X, Y) \pm h(X, Y)] = E[g(X, Y)] \pm E[h(X, Y)].$$

Prueba: Por la definición 4.2,

$$\begin{aligned} E[g(X, Y) \pm h(X, Y)] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} [g(x, y) \pm h(x, y)]f(x, y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y)f(x, y) dx dy \pm \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} h(x, y)f(x, y) dx dy \\ &= E[g(X, Y)] \pm E[h(X, Y)]. \end{aligned} \quad \text{J}$$

Corolario 4.3: Si establecemos que $g(X, Y) = g(X)$ y $h(X, Y) = h(Y)$, vemos que

$$E[g(X) \pm h(Y)] = E[g(X)] \pm E[h(Y)].$$

Corolario 4.4: Si establecemos que $g(X, Y) = X$ y $h(X, Y) = Y$, vemos que

$$E[X \pm Y] = E[X] \pm E[Y].$$

Si X representa la producción diaria de algún artículo de la máquina A y Y la producción diaria del mismo artículo de la máquina B , entonces $X + Y$ representa la cantidad total de artículos que ambas máquinas producen diariamente. El corolario 4.4 establece que la producción diaria promedio para ambas máquinas es igual a la suma de la producción diaria promedio de cada máquina.

Teorema 4.8: Sean X y Y dos variables aleatorias independientes. Entonces,

$$E(XY) = E(X)E(Y).$$

Prueba: Por la definición 4.2,

$$E(XY) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xyf(x, y) dx dy.$$

Como X y Y son independientes, podemos escribir

$$f(x, y) = g(x)h(y),$$

donde $g(x)$ y $h(y)$ son las distribuciones marginales de X y Y , respectivamente. En consecuencia,

$$\begin{aligned} E(XY) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xyg(x)h(y) dx dy = \int_{-\infty}^{\infty} xg(x) dx \int_{-\infty}^{\infty} yh(y) dy \\ &= E(X)E(Y). \end{aligned}$$

Para variables discretas el teorema 4.8 se ilustra mediante un experimento en el que se lanzan un dado verde y uno rojo. La variable aleatoria X representa el resultado de lanzar el dado verde y la variable aleatoria Y el resultado de lanzar el dado rojo. Entonces XY representa el producto de los números que resultan de lanzar el par de dados. A la larga el promedio de los productos de los números es igual al producto del número promedio que resulta de lanzar el dado verde y el número promedio que resulta de lanzar el dado rojo.

Corolario 4.5: Sean X y Y dos variables aleatorias independientes. Entonces, $\sigma_{XY} = 0$.

Prueba: La demostración se puede realizar utilizando los teoremas 4.4 y 4.8.

Ejemplo 4.21: Se sabe que la proporción de galio y arseniuro no afecta el funcionamiento de las obleas de arseniuro de galio que son los principales componentes de los circuitos integrados. Denotemos con X la proporción de galio a arseniuro y con Y el porcentaje de obleas funcionales producidas durante una hora. X y Y son variables aleatorias independientes con la siguiente función de densidad conjunta

$$f(x, y) = \begin{cases} \frac{x(1+3y^2)}{4}, & 0 < x < 2, 0 < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Demuestre que $E(XY) = E(X)E(Y)$, como sugiere el teorema 4.8.

Solución: Por definición,

$$E(XY) = \int_0^1 \int_0^2 \frac{x^2 y (1 + 3y^2)}{4} dx dy = \frac{5}{6}, \quad E(X) = \frac{4}{3}, \quad y \quad E(Y) = \frac{5}{8}.$$

Por lo tanto,

$$E(X)E(Y) = \left(\frac{4}{3}\right) \left(\frac{5}{8}\right) = \frac{5}{6} = E(XY).$$

Concluimos esta sección con la demostración de un teorema y la presentación de varios corolarios que son útiles para calcular varianzas o desviaciones estándar.

Teorema 4.9: Si X y Y son variables aleatorias con distribución de probabilidad conjunta $f(x, y)$, y a , b y c son constantes, entonces

$$\sigma_{aX + bY + c}^2 = a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab\sigma_{XY}.$$

Prueba: Por definición, $\sigma_{aX + bY + c}^2 = E\{[(aX + bY + c) - \mu_{aX + bY + c}]^2\}$. Entonces,

$$\mu_{aX + bY + c} = E(aX + bY + c) = aE(X) + bE(Y) + c = a\mu_X + b\mu_Y + c,$$

si utilizamos el corolario 4.4 y después el corolario 4.2. Por lo tanto,

$$\begin{aligned} \sigma_{aX + bY + c}^2 &= E\{[a(X - \mu_X) + b(Y - \mu_Y)]^2\} \\ &= a^2 E[(X - \mu_X)^2] + b^2 E[(Y - \mu_Y)^2] + 2abE[(X - \mu_X)(Y - \mu_Y)] \\ &= a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab\sigma_{XY}. \end{aligned}$$

Si utilizamos el teorema 4.9, tenemos los siguientes corolarios.

Corolario 4.6: Si se establece que $b = 0$, vemos que

$$\sigma_{aX + c}^2 = a^2 \sigma_X^2 = a^2 \sigma^2.$$

Corolario 4.7: Si se establece que $a = 1$ y $b = 0$, vemos que

$$\sigma_{X + c}^2 = \sigma_X^2 = \sigma^2.$$

Corolario 4.8: Si se establece que $b = 0$ y $c = 0$, vemos que

$$\sigma_{aX}^2 = a^2 \sigma_X^2 = a^2 \sigma^2.$$

Los corolarios 4.6 y 4.7 establecen que la varianza no cambia si se suma o se resta una constante a una variable aleatoria. La suma o resta de una constante simplemente corre los valores de X a la derecha o a la izquierda, pero no cambia su variabilidad. Sin embargo, si una variable aleatoria se multiplica por una constante o se divide entre ésta, entonces los corolarios 4.6 y 4.8 establecen que la varianza se multiplica por el cuadrado de la constante o se divide entre éste.

Corolario 4.9: Si X y Y son variables aleatorias independientes, entonces

$$\sigma_{aX+bY}^2 = a^2\sigma_X^2 + b^2\sigma_Y^2.$$

El resultado que se establece en el corolario 4.9 se obtiene a partir del teorema 4.9 recurriendo al corolario 4.5.

Corolario 4.10: Si X y Y son variables aleatorias independientes, entonces,

$$\sigma_{aX-bY}^2 = a^2\sigma_X^2 + b^2\sigma_Y^2.$$

El corolario 4.10 se obtiene reemplazando b por $-b$ en el corolario 4.9. Al generalizar a una combinación lineal de n variables aleatorias independientes, resulta el corolario 4.11.

Corolario 4.11: Si X_1, X_2, \dots, X_n son variables aleatorias independientes, entonces

$$\sigma_{a_1X_1+a_2X_2+\dots+a_nX_n}^2 = a_1^2\sigma_{X_1}^2 + a_2^2\sigma_{X_2}^2 + \dots + a_n^2\sigma_{X_n}^2.$$

Ejemplo 4.22: Si X y Y son variables aleatorias con varianzas $\sigma_X^2 = 2$ y $\sigma_Y^2 = 4$ y covarianza $\sigma_{XY} = -2$, calcule la varianza de la variable aleatoria $Z = 3X - 4Y + 8$.

Solución:

$$\begin{aligned}\sigma_Z^2 &= \sigma_{3X-4Y+8}^2 = \sigma_{3X-4Y}^2 && \text{(por el corolario 4.6)} \\ &= 9\sigma_X^2 + 16\sigma_Y^2 - 24\sigma_{XY} && \text{(por el teorema 4.9)} \\ &= (9)(2) + (16)(4) - (24)(-2) = 130. && \blacksquare\end{aligned}$$

Ejemplo 4.23: Denotemos con X y Y la cantidad de dos tipos diferentes de impurezas en un lote de cierto producto químico. Suponga que X y Y son variables aleatorias independientes con varianzas $\sigma_X^2 = 2$ y $\sigma_Y^2 = 3$. Calcule la varianza de la variable aleatoria $Z = 3X - 2Y + 5$.

Solución:

$$\begin{aligned}\sigma_Z^2 &= \sigma_{3X-2Y+5}^2 = \sigma_{3X-2Y}^2 && \text{(por el corolario 4.6)} \\ &= 9\sigma_X^2 + 4\sigma_Y^2 && \text{(por el corolario 4.10)} \\ &= (9)(2) + (4)(3) = 30. && \blacksquare\end{aligned}$$

¿Qué sucede si la función es no lineal?

En las secciones anteriores estudiamos propiedades de funciones lineales de variables aleatorias por razones muy importantes. En los capítulos 8 a 15 se estudiarán y ejemplificarán problemas de la vida real, en los cuales el analista construye un **modelo lineal** para describir un conjunto de datos y, en consecuencia, describir o explicar el comportamiento de un fenómeno científico. Así que resulta natural que encontremos los valores esperados y las varianzas de combinaciones lineales de variables aleatorias. Sin embargo, hay situaciones en que las propiedades de las funciones **no lineales** de variables aleatorias se vuelven importantes. En efecto, hay muchos fenómenos científicos de naturaleza no lineal, donde el modelado estadístico que utiliza funciones no lineales adquiere gran importancia. De hecho, en el capítulo 12 se estudia el modelado de los que se han convertido en modelos estándar no lineales. En realidad, incluso una función simple de variables aleatorias, como $Z = XY$, ocurre con bastante frecuencia en la prác-

tica, y a diferencia del caso del valor esperado de las combinaciones lineales de variables aleatorias, no hay una simple regla general. Por ejemplo,

$$E(Z) = E(X/Y) \neq E(X)/E(Y),$$

excepto en circunstancias muy especiales.

El material dado por los teoremas 4.5 a 4.9 y los diversos corolarios son sumamente útiles, ya que no hay restricciones sobre la forma de la densidad o las funciones de probabilidad, aparte de la propiedad de independencia cuando ésta se requiere, como en los corolarios posteriores al teorema 4.9. Para ilustrar considere el ejemplo 4.23; la varianza de $Z = 3X - 2Y + 5$ no requiere restricciones en las distribuciones de las cantidades X y Y de los dos tipos de impurezas. Sólo se requiere la independencia entre X y Y . Por consiguiente, disponemos de la capacidad de calcular $\mu_{g(X)}$ y $\sigma_{g(X)}^2$ para cualquier función $g(\cdot)$ a partir de los principios iniciales establecidos en los teoremas 4.1 y 4.3, donde se supone que se conoce la distribución $f(x)$ correspondiente. Los ejercicios 4.40, 4.41 y 4.42, entre otros, ilustran el uso de tales teoremas. De modo que, si $g(x)$ es una función no lineal y se conoce la función de densidad (o función de probabilidad en el caso discreto), $\mu_{g(X)}$ y $\sigma_{g(X)}^2$ pueden evaluarse con exactitud. No obstante, como en el caso de las reglas dadas para combinaciones lineales, ¿habría reglas para funciones no lineales que se puedan utilizar cuando no se conoce la forma de la distribución de las variables aleatorias pertinentes?

En general, suponga que X es una variable aleatoria y que $Y = g(x)$. La solución general para $E(Y)$ o $\text{Var}(Y)$ puede ser difícil y depende de la complejidad de la función $g(\cdot)$. Sin embargo, hay aproximaciones disponibles que dependen de una aproximación lineal de la función $g(x)$. Por ejemplo, suponga que denotamos $E(X)$ como μ y $\text{Var}(X) = \sigma_X^2$. Entonces, una aproximación a las series de Taylor de $g(x)$ alrededor de $X = \mu_X$ da

$$g(x) = g(\mu_X) + \left. \frac{\partial g(x)}{\partial x} \right|_{x=\mu_X} (x - \mu_X) + \left. \frac{\partial^2 g(x)}{\partial x^2} \right|_{x=\mu_X} \frac{(x - \mu_X)^2}{2} + \dots$$

Como resultado, si truncamos después el término lineal y tomamos el valor esperado de ambos lados, obtenemos $E[g(X)] \approx g(\mu_X)$, que ciertamente es intuitivo y en algunos casos ofrece una aproximación razonable. No obstante, si incluimos el término de segundo orden de la serie de Taylor, entonces tenemos un ajuste de segundo orden para esta *aproximación de primer orden* como sigue:

Aproximación de
 $E[g(X)]$

$$E[g(X)] \approx g(\mu_X) + \left. \frac{\partial^2 g(x)}{\partial x^2} \right|_{x=\mu_X} \frac{\sigma_X^2}{2}.$$

Ejemplo 4.24: Dada la variable aleatoria X con media μ_X y varianza σ_X^2 , determine la aproximación de segundo orden para $E(e^X)$.

Solución: Como $\frac{\partial e^x}{\partial x} = e^x$ y $\frac{\partial^2 e^x}{\partial x^2} = e^x$, obtenemos $E(e^X) \approx e^{\mu_X} (1 + \sigma_X^2/2)$. ▮

De manera similar, podemos desarrollar una aproximación para $\text{Var}[g(x)]$ tomando la varianza de ambos lados de la expansión de la serie de Taylor de primer orden de $g(x)$.

Aproximación de
 $\text{Var}[g(X)]$

$$\text{Var}[g(X)] \approx \left[\left. \frac{\partial g(x)}{\partial x} \right|_{x=\mu_X} \right]^2 \sigma_X^2.$$

Ejemplo 4.25: Dada la variable aleatoria X , como en el ejemplo 4.24, determine una fórmula aproximada para $\text{Var}[g(x)]$.

Solución: De nuevo, $\frac{\partial e^x}{\partial x} = e^x$ por lo tanto, $\text{Var}(X) \approx e^{2\mu_x} \sigma_x^2$. ▀

Estas aproximaciones se pueden extender a las funciones no lineales de más de una variable aleatoria.

Dado un conjunto de variables aleatorias independientes X_1, X_2, \dots, X_k con medias $\mu_1, \mu_2, \dots, \mu_k$ y varianzas $\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2$, respectivamente, sea

$$Y = h(X_1, X_2, \dots, X_k)$$

una función no lineal; entonces tenemos las siguientes aproximaciones para $E(Y)$ y $\text{Var}(Y)$:

$$E(Y) \approx h(\mu_1, \mu_2, \dots, \mu_k) + \sum_{i=1}^k \frac{\sigma_i^2}{2} \left[\frac{\partial^2 h(x_1, x_2, \dots, x_k)}{\partial x_i^2} \right] \Bigg|_{x_i = \mu_i, 1 \leq i \leq k},$$

$$\text{Var}(Y) \approx \sum_{i=1}^k \left[\frac{\partial h(x_1, x_2, \dots, x_k)}{\partial x_i} \right]^2 \Bigg|_{x_i = \mu_i, 1 \leq i \leq k} \sigma_i^2.$$

Ejemplo 4.26: Considere dos variables aleatorias independientes X y Z , con medias μ_x, μ_z y varianzas σ_x^2 y σ_z^2 , respectivamente. Considere una variable aleatoria

$$Y = X/Z.$$

Determine aproximaciones para $E(Y)$ y $\text{Var}(Y)$.

Solución: Para $E(Y)$, debemos usar $\frac{\partial y}{\partial x} = \frac{1}{z}$ y $\frac{\partial y}{\partial z} = -\frac{x}{z^2}$. Por consiguiente,

$$\frac{\partial^2 y}{\partial x^2} = 0 \quad \text{y} \quad \frac{\partial^2 y}{\partial z^2} = \frac{2x}{z^3}.$$

Como resultado,

$$E(Y) \approx \frac{\mu_x}{\mu_z} + \frac{\mu_x}{\mu_z^3} \sigma_z^2 = \frac{\mu_x}{\mu_z} \left(1 + \frac{\sigma_z^2}{\mu_z^2} \right),$$

y la aproximación para la varianza de Y está dada por

$$\text{Var}(Y) \approx \frac{1}{\mu_z^2} \sigma_x^2 + \frac{\mu_x^2}{\mu_z^4} \sigma_z^2 = \frac{1}{\mu_z^2} \left(\sigma_x^2 + \frac{\mu_x^2}{\mu_z^2} \sigma_z^2 \right). \quad \blacksquare$$

4.4 Teorema de Chebyshev

En la sección 4.2 establecimos que la varianza de una variable aleatoria nos dice algo acerca de la variabilidad de las observaciones con respecto a la media. Si una variable aleatoria tiene una varianza o desviación estándar pequeña, esperaríamos que la mayoría de los valores se agrupen alrededor de la media. Por lo tanto, la probabilidad de que una variable aleatoria tome un valor dentro de cierto intervalo alrededor de la media es mayor que para una variable aleatoria similar con una desviación estándar mayor. Si pensamos en la probabilidad en términos de área, esperaríamos una distribución continua con un valor grande de σ para indicar una variabilidad mayor y, por lo tanto, esperaríamos que el área esté más extendida, como en la figura 4.2(a). Una distribución con una desviación estándar pequeña debería tener la mayor parte de su área cercana a μ , como en la figura 4.2(b).

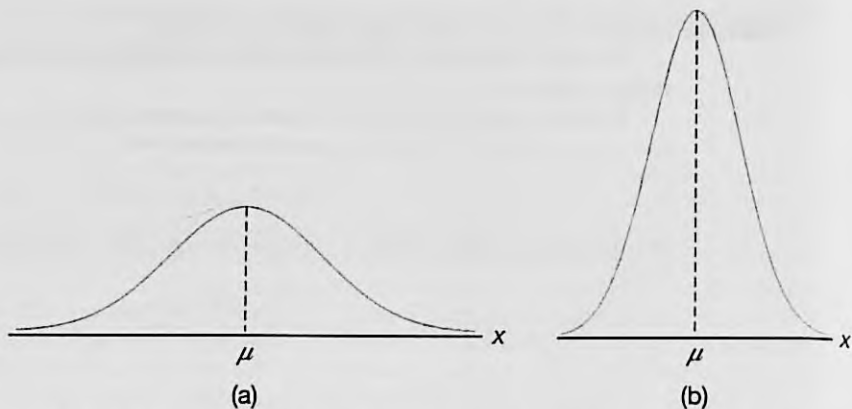


Figura 4.2: Variabilidad de observaciones continuas alrededor de la media.

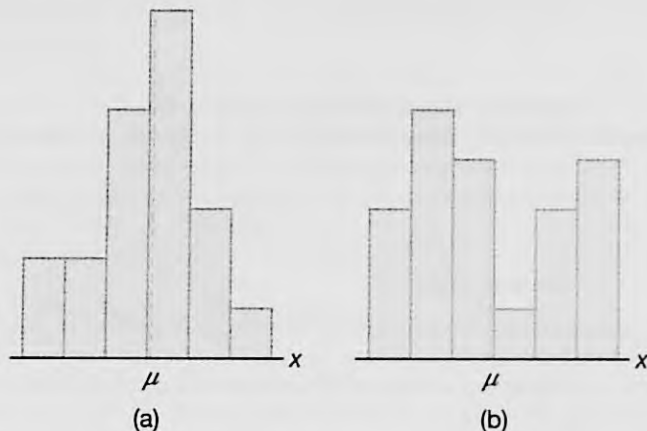


Figura 4.3: Variabilidad de observaciones discretas alrededor de la media.

Podemos argumentar lo mismo para una distribución discreta. En el histograma de probabilidad de la figura 4.3(b) el área se extiende mucho más que en la figura 4.3(a), lo cual indica una distribución más variable de mediciones o resultados.

El matemático ruso P. L. Chebyshev (1821-1894) descubrió que la fracción del área entre cualesquiera dos valores simétricos alrededor de la media está relacionada con la desviación estándar. Como el área bajo una curva de distribución de probabilidad, o la de un histograma de probabilidad, suma 1, el área entre cualesquiera dos números es la probabilidad de que la variable aleatoria tome un valor entre estos números.

El siguiente teorema, planteado por Chebyshev, ofrece una estimación conservadora de la probabilidad de que una variable aleatoria tome un valor dentro de k desviaciones estándar de su media para cualquier número real k .

Teorema 4.10: (Teorema de Chebyshev) La probabilidad de que cualquier variable aleatoria X tome un valor dentro de k desviaciones estándar de la media es de al menos $1 - 1/k^2$. Es decir,

$$P(\mu - k\sigma < X < \mu + k\sigma) \geq 1 - \frac{1}{k^2}.$$

Para $k = 2$ el teorema establece que la variable aleatoria X tiene una probabilidad de al menos $1 - 1/2^2 = 3/4$ de caer dentro de dos desviaciones estándar a partir de la media; es decir, que tres cuartas partes o más de las observaciones de cualquier distribución se localizan en el intervalo $\mu \pm 2\sigma$. De manera similar, el teorema afirma que al menos ocho novenos de las observaciones de cualquier distribución caen en el intervalo $\mu \pm 3\sigma$.

Ejemplo 4.27: Una variable aleatoria X tiene una media $\mu = 8$, una varianza $\sigma^2 = 9$ y una distribución de probabilidad desconocida. Calcule

a) $P(-4 < X < 20)$,

b) $P(|X - 8| \geq 6)$.

Solución: a) $P(-4 < X < 20) = P[8 - (4)(3) < X < 8 + (4)(3)] \geq \frac{15}{16}$.

b) $P(|X - 8| \geq 6) = 1 - P(|X - 8| < 6) = 1 - P(-6 < X - 8 < 6)$
 $= 1 - P[8 - (2)(3) < X < 8 + (2)(3)] \leq \frac{1}{4}$. ▮

El teorema de Chebyshev tiene validez para cualquier distribución de observaciones, por lo cual los resultados generalmente son débiles. El valor que proporciona el teorema es sólo un límite inferior, es decir, sabemos que la probabilidad de una variable aleatoria que cae dentro de dos desviaciones estándar de la media *no puede ser menor* que $3/4$, pero nunca sabemos cuánto podría ser en realidad. Sólo cuando conocemos la distribución de probabilidad podemos determinar probabilidades exactas. Por esta razón llamamos al teorema resultado de *distribución libre*. Cuando se supongan distribuciones específicas, como ocurrirá en los siguientes capítulos, los resultados serán menos conservadores. El uso del teorema de Chebyshev se restringe a situaciones donde se desconoce la forma de la distribución.

Ejercicios

4.53 Remítase al ejercicio 4.35 de la página 127 y calcule la media y la varianza de la variable aleatoria discreta $Z = 3X - 2$, donde X representa el número de errores por 100 líneas de código.

4.54 Use el teorema 4.5 y el corolario 4.6 para calcular la media y la varianza de la variable aleatoria $Z = 5X + 3$, donde X tiene la distribución de probabilidad del ejercicio 4.36 de la página 127.

4.55 Suponga que una tienda de abarrotes compra 5 envases de leche descremada al precio de mayoreo de \$1.20 por envase y la vende a \$1.65 por envase. Después de la fecha de caducidad, la leche que no se vende se retira de los anaqueles y el tendero recibe un crédito del distribuidor igual a tres cuartas partes del

precio de mayoreo. Si la distribución de probabilidad de la variable aleatoria es X y el número de envases que se venden de este lote es

x	0	1	2	3	4	5
$f(x)$	$\frac{1}{15}$	$\frac{2}{15}$	$\frac{2}{15}$	$\frac{3}{15}$	$\frac{4}{15}$	$\frac{5}{15}$

calcule la utilidad esperada.

4.56 Repita el ejercicio 4.43 de la página 127 aplicando el teorema 4.5 y el corolario 4.6.

4.57 Sea X una variable aleatoria con la siguiente distribución de probabilidad:

x	-3	6	9
$f(x)$	$\frac{1}{6}$	$\frac{1}{2}$	$\frac{1}{3}$

Calcule $E(X)$ y $E(X^2)$ y luego utilice estos valores para evaluar $E[(2X + 1)^2]$.

4.58 El tiempo total que una adolescente utiliza su secadora de pelo durante un año, medido en unidades de 100 horas, es una variable aleatoria continua X que tiene la siguiente función de densidad

$$f(x) = \begin{cases} x, & 0 < x < 1, \\ 2 - x, & 1 \leq x < 2, \\ 0, & \text{en otro caso.} \end{cases}$$

Utilice el teorema 4.6 para evaluar la media de la variable aleatoria $Y = 60X^2 + 39X$, donde Y es igual al número de kilowatts-hora que gasta al año.

4.59 Si una variable aleatoria X se define de manera que

$$E[(X - 1)^2] = 10 \quad \text{y} \quad E[(X - 2)^2] = 6,$$

calcule μ y σ^2 .

4.60 Suponga que X y Y son variables aleatorias independientes que tienen la siguiente distribución de probabilidad conjunta

$f(x, y)$		x	
		2	4
y	1	0.10	0.15
	3	0.20	0.30
	5	0.10	0.15

Calcule

- a) $E(2X - 3Y)$;
b) $E(XY)$.

4.61 Use el teorema 4.7 para evaluar $E(2XY^2 - X^2Y)$ en la distribución de probabilidad conjunta que se muestra en la tabla 3.1 de la página 96.

4.62 Si X y Y son variables aleatorias independientes con varianzas $\sigma_X^2 = 5$ y $\sigma_Y^2 = 3$, calcule la varianza de la variable aleatoria $Z = -2X + 4Y - 3$.

4.63 Repita el ejercicio 4.62 si X y Y no son independientes y $\sigma_{XY} = 1$

4.64 Suponga que X y Y son variables aleatorias independientes con densidades de probabilidad y

$$g(x) = \begin{cases} \frac{8}{x^3}, & x > 2, \\ 0, & \text{en otro caso,} \end{cases}$$

y

$$h(y) = \begin{cases} 2y, & 0 < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el valor esperado de $Z = XY$.

4.65 Sea X el número que resulta cuando se lanza un dado rojo y Y el número que resulta cuando se lanza un dado verde. Calcule

- a) $E(X + Y)$;
b) $E(X - Y)$;
c) $E(XY)$.

4.66 Sea X el número que resulta cuando se lanza un dado verde y Y el número que resulta cuando se lanza un dado rojo. Calcule la varianza de la variable aleatoria

- a) $2X - Y$;
b) $X + 3Y - 5$.

4.67 Si la función de densidad conjunta de X y Y está dada por

$$f(x, y) = \begin{cases} \frac{2}{7}(x + 2y), & 0 < x < 1, 1 < y < 2, \\ 0, & \text{en otro caso,} \end{cases}$$

calcule el valor esperado de $g(X, Y) = \frac{X}{Y^3} + X^2Y$.

4.68 Se sabe que la potencia P en watts que se disipa en un circuito eléctrico con resistencia R está dada por $P = I^2R$, donde I es la corriente en amperes y R es una constante fija en 50 ohms. Sin embargo, I es una variable aleatoria con $\mu_I = 15$ amperes y $\sigma_I^2 = 0.03$ amperes². Dé aproximaciones numéricas a la media y a la varianza de la potencia P .

4.69 Considere el ejercicio de repaso 3.77 de la página 108. Las variables aleatorias X y Y representan el número de vehículos que llegan a dos esquinas de calles separadas durante cierto periodo de 2 minutos en el día. La distribución conjunta es

$$f(x, y) = \left(\frac{1}{4^{(x+y)}} \right) \left(\frac{9}{16} \right),$$

para $x = 0, 1, 2, \dots$, y $y = 0, 1, 2, \dots$

- a) Determine $E(X)$, $E(Y)$, $\text{Var}(X)$ y $\text{Var}(Y)$.
b) Considere que $Z = X + Y$ es la suma de ambas. Calcule $E(Z)$ y $\text{Var}(Z)$.

4.70 Considere el ejercicio de repaso 3.64 de la página 107. Hay dos líneas de servicio. Las variables aleatorias X y Y son las proporciones del tiempo que la línea 1 y la línea 2 están en funcionamiento, respectivamente. La función de densidad de probabilidad conjunta para (X, Y) está dada por

$$f(x, y) = \begin{cases} \frac{3}{2}(x^2 + y^2), & 0 \leq x, y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

a) Determine si X y Y son independientes o no.

- b) Se tiene interés por saber algo acerca de la proporción de $Z = X + Y$, la suma de las dos proporciones. Calcule $E(X + Y)$. También calcule $E(XY)$.
- c) Calcule $\text{Var}(X)$, $\text{Var}(Y)$ y $\text{Cov}(X, Y)$.
- d) Calcule $\text{Var}(X + Y)$.

4.71 El periodo Y en minutos que se requiere para generar un reflejo humano ante el gas lacrimógeno tiene la siguiente función de densidad

$$f(y) = \begin{cases} \frac{1}{4}e^{-y/4}, & 0 \leq y < \infty, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) ¿Cuál es el tiempo medio para el reflejo?
- b) Calcule $E(Y^2)$ y $\text{Var}(Y)$.

4.72 Una empresa industrial desarrolló una máquina de limpiar alfombras con buen rendimiento de combustible porque limpia más superficie de alfombra en menos tiempo. Se tiene interés por una variable aleatoria Y , la cantidad en galones por minuto que ofrece. Se sabe que la función de densidad está dada por

$$f(y) = \begin{cases} 1, & 7 \leq y \leq 8, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Determine la función de densidad.
- b) Calcule $E(Y)$, $E(Y^2)$ y $\text{Var}(Y)$.

4.73 Para la situación del ejercicio 4.72 calcule $E(e^Y)$ utilizando el teorema 4.1, es decir, mediante el uso de

$$E(e^Y) = \int_7^8 e^y f(y) dy.$$

Luego, calcule $E(e^Y)$ sin utilizar $f(y)$. En su lugar utilice el ajuste de segundo orden para la aproximación de primer orden de $E(e^Y)$. Comente al respecto.

4.74 Considere nuevamente la situación del ejercicio 4.72, donde se le pide calcular $\text{Var}(e^Y)$. Utilice los teoremas 4.2 y 4.3 y defina $Z = e^Y$. En consecuencia, utilice las condiciones del ejercicio 4.73 para calcular

$$\text{Var}(Z) = E(Z^2) - [E(Z)]^2.$$

Luego hágalo sin utilizar $f(y)$. En su lugar utilice la aproximación de primer orden a las series de Taylor para $\text{Var}(e^Y)$. ¡Comente al respecto!

4.75 Una empresa eléctrica fabrica una bombilla de luz de 100 watts que, de acuerdo con las especificaciones escritas en la caja, tiene una vida media de 900 horas con una desviación estándar de 50 horas. A lo sumo, ¿qué porcentaje de las bombillas no duran al menos 700 horas? Suponga que la distribución es simétrica alrededor de la media.

4.76 En una planta de ensamble automotriz se crean 70 nuevos puestos de trabajo y se presentan 1000 aspirantes. Para seleccionar entre los aspirantes a los 70 mejores la armadora aplica un examen que abarca habilidad mecánica, destreza manual y capacidad matemática. La calificación media de este examen resulta ser 60 y las calificaciones tienen una desviación estándar de 6. ¿Una persona que obtiene una calificación de 84 puede obtener uno de los puestos? [Sugerencia: Utilice el teorema de Chebyshev]. Suponga que la distribución es simétrica alrededor de la media.

4.77 Una variable aleatoria X tiene una media $\mu = 10$ y una varianza $\sigma^2 = 4$. Utilice el teorema de Chebyshev para calcular

- a) $P(|X - 10| \geq 3)$;
- b) $P(|X - 10| < 3)$;
- c) $P(5 < X < 15)$;

d) el valor de la constante c tal que

$$P(|X - 10| \geq c) \leq 0.04.$$

4.78 Calcule $P(\mu - 2\sigma < X < \mu + 2\sigma)$, donde X tiene la siguiente función de densidad

$$f(x) = \begin{cases} 6x(1-x), & 0 < x < 1, \\ 0, & \text{en otro caso,} \end{cases}$$

y compare con el resultado dado por el teorema de Chebyshev.

Ejercicios de repaso

4.79 Demuestre el teorema de Chebyshev.

4.80 Calcule la covarianza de las variables aleatorias X y Y que tienen la siguiente función de densidad de probabilidad conjunta

$$f(x, y) = \begin{cases} x + y, & 0 < x < 1, 0 < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

4.81 Remítase a las variables aleatorias cuya función de densidad de probabilidad conjunta está dada en el ejercicio 3.47 de la página 105 y calcule la cantidad promedio de queroseno que queda en el tanque al final del día.

4.82 Suponga que la duración X en minutos de un tipo específico de conversación telefónica es una variable aleatoria con función de densidad de probabilidad

$$f(x) = \begin{cases} \frac{1}{5}e^{-x/5}, & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Determine la duración media $E(X)$ de este tipo de conversación telefónica.
 b) Calcule la varianza y la desviación estándar de X .
 c) Calcule $E[(X + 5)^2]$.

4.83 Remítase a las variables aleatorias cuya función de densidad conjunta está dada en el ejercicio 3.41 de la página 105 y calcule la covarianza entre el peso de las cremas y el peso de los chiclosos en estas cajas de chocolates.

4.84 Remítase a las variables aleatorias cuya función de densidad de probabilidad conjunta está dada en el ejercicio 3.41 de la página 105 y calcule el peso esperado para la suma de las cremas y los chiclosos si uno compra una caja de tales chocolates.

4.85 Suponga que se sabe que la vida de un compresor particular X , en horas, tiene la siguiente función de densidad

$$f(x) = \begin{cases} \frac{1}{900}e^{-x/900}, & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Calcule la vida media del compresor.
 b) Calcule $E(X^2)$.
 c) Calcule la varianza y la desviación estándar de la variable aleatoria X .

4.86 Remítase a las variables aleatorias cuya función de densidad conjunta está dada en el ejercicio 3.40 de la página 105,

- a) calcule μ_x y μ_y ;
 b) calcule $E[(X + Y)/2]$.

4.87 Demuestre que $\text{Cov}(aX, bY) = ab \text{Cov}(X, Y)$.

4.88 Considere la función de densidad del ejercicio de repaso 4.85. Demuestre que el teorema de Chebyshev es válido para $k = 2$ y $k = 3$.

4.89 Considere la siguiente función de densidad conjunta

$$f(x, y) = \begin{cases} \frac{16xy}{x^3}, & x > 2, 0 < y < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Calcule el coeficiente de correlación ρ_{XY} .

4.90 Considere las variables aleatorias X y Y del ejercicio 4.63 de la página 138. Calcule ρ_{XY} .

4.91 La utilidad de un distribuidor, en unidades de \$5000, por un automóvil nuevo es una variable aleatoria X que tiene la siguiente función de densidad

$$f(x) = \begin{cases} 2(1-x), & 0 \leq x \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Calcule la varianza de la utilidad del distribuidor.
 b) Demuestre que el teorema de Chebyshev es válido para $k = 2$ con la función de densidad anterior.

c) ¿Cuál es la probabilidad de que la utilidad exceda \$500?

4.92 Considere el ejercicio 4.10 de la página 117. ¿Se puede decir que las calificaciones dadas por los dos expertos son independientes? Explique su respuesta.

4.93 Los departamentos de marketing y de contabilidad de una empresa determinaron que si la empresa comercializa su producto creado recientemente, su contribución a las utilidades de la empresa durante los próximos 6 meses será la siguiente:

Contribución a las utilidades	Probabilidad
-\$5,000	0.2
\$10,000	0.5
\$30,000	0.3

¿Cuál es la utilidad esperada de la empresa?

4.94 En un sistema de apoyo para el programa espacial estadounidense, un componente crucial único funciona sólo 85 por ciento del tiempo. Para aumentar la confiabilidad del sistema se decidió instalar tres componentes paralelos, de manera que el sistema falle sólo si todos fallan. Suponga que los componentes actúan de forma independiente y que son equivalentes en el sentido de que 3 de ellos tienen una tasa de éxito de 85 por ciento. Considere la variable aleatoria X como el número de componentes de cada tres que fallan.

- a) Escriba una función de probabilidad para la variable aleatoria X .
 b) ¿Cuál es $E(X)$ (es decir, el número medio de componentes de cada tres que fallan)?
 c) ¿Cuál es $\text{Var}(X)$?
 d) ¿Cuál es la probabilidad de que el sistema completo sea exitoso?
 e) ¿Cuál es la probabilidad de que falle el sistema?
 f) Si se desea que el sistema tenga una probabilidad de éxito de 0.99, ¿son suficientes los tres componentes? Si no lo son, ¿cuántos se requerirían?

4.95 En los negocios es importante planear y llevar a cabo investigación para anticipar lo que ocurrirá al final del año. La investigación sugiere que el espectro de utilidades (pérdidas) de cierta empresa, con sus respectivas probabilidades, es el siguiente:

Utilidad	Probabilidad
-\$15,000	0.05
\$0	0.15
\$15,000	0.15
\$25,000	0.30
\$40,000	0.15
\$50,000	0.10
\$100,000	0.05
\$150,000	0.03
\$200,000	0.02

- a) ¿Cuál es la utilidad esperada?
 b) Determine la desviación estándar de las utilidades.

4.96 Mediante un conjunto de datos, y por la amplia investigación, se sabe que la cantidad de tiempo que cierto empleado de una empresa llega tarde a trabajar, medido en segundos, es una variable aleatoria X con la siguiente función de densidad

$$f(x) = \begin{cases} \frac{3}{(4)(50^3)}(50^2 - x^2), & -50 \leq x \leq 50, \\ 0, & \text{en otro caso.} \end{cases}$$

En otras palabras, él no sólo llega ligeramente retrasado a veces, sino que también puede llegar temprano a trabajar.

- a) Calcule el valor esperado del tiempo en segundos que llega tarde.
 b) Calcule $E(X^2)$.
 c) ¿Cuál es la desviación estándar del tiempo en que llega tarde?

4.97 Un camión de carga viaja desde el punto A hasta el punto B y regresa por la misma ruta diariamente. Hay cuatro semáforos en la ruta. Sea X_1 el número de semáforos en rojo que el camión encuentra cuando va de A a B y X_2 el número de los que encuentra en el viaje de regreso. Los datos recabados durante un periodo largo sugieren que la distribución de probabilidad conjunta para (X_1, X_2) está dada por

x_1	x_2				
	0	1	2	3	4
0	0.01	0.01	0.03	0.07	0.01
1	0.03	0.05	0.08	0.03	0.02
2	0.03	0.11	0.15	0.01	0.01
3	0.02	0.07	0.10	0.03	0.01
4	0.01	0.06	0.03	0.01	0.01

- a) Determine la densidad marginal de X_1 .
 b) Determine la densidad marginal de X_2 .
 c) Determine la distribución de densidad condicional de X_1 dado que $X_2 = 3$.
 d) Determine $E(X_1)$.
 e) Determine $E(X_2)$.
 f) Determine $E(X_1 | X_2 = 3)$.
 g) Determine la desviación estándar de X_1 .

4.98 Una tienda de abarrotes tiene dos sitios separados en sus instalaciones donde los clientes pueden pagar cuando se marchan. Estos dos lugares tienen dos cajas registradoras y dos empleados que atienden a los clientes que van a pagar. Sea X el número de la caja registradora que se utiliza en un momento específico en el sitio 1 y Y el número de la caja registradora que se utiliza en el mismo momento en el sitio 2. La función de probabilidad conjunta está dada por

x	y		
	0	1	2
0	0.12	0.04	0.04
1	0.08	0.19	0.05
2	0.06	0.12	0.30

- a) Determine la densidad marginal de X y de Y , así como la distribución de probabilidad de X , dado que $Y = 2$.
 b) Determine $E(X)$ y $\text{Var}(X)$.
 c) Determine $E(X|Y = 2)$ y $\text{Var}(X|Y = 2)$.

4.99 Considere un transbordador que puede llevar tanto autobuses como automóviles en un recorrido a través de una vía fluvial. Cada viaje cuesta al propietario aproximadamente \$10. La tarifa por automóvil es de \$3 y por autobús es de \$8. Sean X y Y el número de autobuses y automóviles, respectivamente, que se transportan en un viaje específico. La distribución conjunta de X y Y está dada por

y	x		
	0	1	2
0	0.01	0.01	0.03
1	0.03	0.08	0.07
2	0.03	0.06	0.06
3	0.07	0.07	0.13
4	0.12	0.04	0.03
5	0.08	0.06	0.02

Calcule la utilidad esperada para el viaje del transbordador.

4.100 Como veremos en el capítulo 12, los métodos estadísticos asociados con los modelos lineal y no lineal son muy importantes. De hecho, a menudo las funciones exponenciales se utilizan en una amplia gama de problemas científicos y de ingeniería. Considere un modelo que se ajusta a un conjunto de datos que implica los valores medidos k_1 y k_2 , y una respuesta específica Y a las mediciones. El modelo postulado es

$$\hat{Y} = e^{b_0 + b_1 k_1 + b_2 k_2},$$

donde \hat{Y} denota el valor estimado de Y , k_1 y k_2 son valores fijos y b_0 , b_1 y b_2 son estimados de constantes y, por lo tanto, variables aleatorias. Suponga que tales variables aleatorias son independientes y use la fórmula aproximada para la varianza de una función no lineal de más de una variable. Dé una expresión para $\text{Var}(\hat{Y})$. Suponga que se conocen las medias de b_0 , b_1 y b_2 y que son β_0 , β_1 y β_2 , y también suponga que se conocen las varianzas de b_0 , b_1 y b_2 y que son σ_0^2 , σ_1^2 y σ_2^2 .

4.101 Considere el ejercicio de repaso 3.73 de la página 108, el cual implica Y , la proporción de impurezas en un lote, donde la función de densidad está dada por

$$f(y) = \begin{cases} 10(1-y)^9, & 0 \leq y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- a) Calcule el porcentaje esperado de impurezas.
 b) Calcule el valor esperado de la proporción de la calidad del material (es decir, calcule $E(1-Y)$).

c) Calcule la varianza de la variable aleatoria $Z = 1 - Y$.

4.102 Proyecto: Sea $X =$ número de horas que cada estudiante del grupo durmió la noche anterior. Cree una variable discreta utilizando los siguientes intervalos arbitrarios:

$X < 3, 3 \leq X < 6, 6 \leq X < 9$ y $X \geq 9$.

- a) Estime la distribución de probabilidad para X .
 b) Calcule la media estimada y la varianza para X .

4.5 Posibles riesgos y errores conceptuales; relación con el material de otros capítulos

El material que se cubrió en este capítulo es fundamental, como el contenido del capítulo 3. Mientras que en el capítulo 3 nos concentramos en las características generales de una distribución de probabilidad, en el presente capítulo definimos cantidades importantes o *parámetros* que caracterizan la naturaleza general del sistema. La **media** de una distribución refleja una *tendencia central*, en tanto que la **varianza** o la **desviación estándar** reflejan *variabilidad* en el sistema. Además, la covarianza refleja la tendencia de dos variables aleatorias a “moverse juntas” en un sistema. Estos importantes parámetros serán fundamentales en el estudio de los siguientes capítulos.

El lector debería comprender que el tipo de distribución a menudo está determinado por el contexto científico. Sin embargo, los valores del parámetro necesitan estimarse a partir de datos científicos. Por ejemplo, en el caso del ejercicio de repaso 4.85 el fabricante del compresor podría saber (material que se presentará en el capítulo 6), por su experiencia y conocimiento del tipo de compresor, que la naturaleza de la distribución es como se indica en el ejercicio. Pero la media $\mu = 900$ se **estimaría** a partir de la experimentación con la máquina. Aunque aquí se da por conocido el valor del parámetro de 900, en situaciones reales eso no ocurrirá sin el uso de datos experimentales. El capítulo 9 se dedica a la **estimación**.

Capítulo 5

Algunas distribuciones de probabilidad discreta

5.1 Introducción y motivación

La distribución de probabilidad discreta describe el comportamiento de una variable aleatoria, independientemente de si se representa de forma gráfica o mediante un histograma, en forma tabular o con una fórmula. A menudo las observaciones que se generan mediante diferentes experimentos estadísticos tienen el mismo tipo general de comportamiento. En consecuencia, las variables aleatorias discretas asociadas con estos experimentos se pueden describir esencialmente con la misma distribución de probabilidad y, por lo tanto, es posible representarlas usando una sola fórmula. De hecho, se necesitan sólo unas cuantas distribuciones de probabilidad importantes para describir muchas de las variables aleatorias discretas que se encuentran en la práctica.

Este conjunto de distribuciones en realidad describe varios fenómenos aleatorios de la vida real. Por ejemplo, en un estudio en el que se probó la eficacia de un nuevo fármaco, de todos los pacientes que lo utilizaron, el número de pacientes que se curaron se aproximó a una distribución binomial (sección 5.2). En un ejemplo en una industria, cuando se prueba una muestra de artículos seleccionados de un lote de producción, el número de productos defectuosos en la muestra por lo general se puede representar como una variable aleatoria hipergeométrica (sección 5.3). En un problema estadístico de control de calidad el experimentador señalará un cambio en la media del proceso cuando los datos observacionales excedan ciertos límites. El número de muestras requeridas para generar una falsa alarma sigue una distribución geométrica, que es un caso especial de distribución binomial negativa (sección 5.4). Por otro lado, el número de leucocitos de una cantidad fija de una muestra de la sangre de un individuo suele ser aleatorio y podría describirse mediante una distribución de Poisson (sección 5.5). En este capítulo se presentarán esas distribuciones de uso común con varios ejemplos.

5.2 Distribuciones binomial y multinomial

Con frecuencia un experimento consta de pruebas repetidas, cada una con dos resultados posibles que se pueden denominar **éxito** o **fracaso**. La aplicación más evidente tiene que ver con la prueba de artículos a medida que salen de una línea de ensamble, donde cada

prueba o experimento puede indicar si un artículo está o no defectuoso. Podemos elegir definir cualquiera de los resultados como éxito. El proceso se conoce como **proceso de Bernoulli** y cada ensayo se denomina **experimento de Bernoulli**. Por ejemplo, si extraemos cartas de una baraja y éstas no se reemplazan, cambian las probabilidades en la repetición de cada ensayo; es decir, la probabilidad de seleccionar una carta de corazones en la primera extracción es $1/4$, pero en la segunda es una probabilidad condicional que tiene un valor de $13/51$ o $12/51$, dependiendo de si resulta un corazón en la primera extracción; entonces éste ya no sería considerado un conjunto de experimentos de Bernoulli.

El proceso de Bernoulli

En términos estrictos el proceso de Bernoulli se caracteriza por lo siguiente:

1. El experimento consta de ensayos repetidos.
2. Cada ensayo produce un resultado que se puede clasificar como éxito o fracaso.
3. La probabilidad de un éxito, que se denota con p , permanece constante de un ensayo a otro.
4. Los ensayos repetidos son independientes.

Considere el conjunto de experimentos de Bernoulli en el que se seleccionan tres artículos al azar de un proceso de producción, luego se inspeccionan y se clasifican como defectuosos o no defectuosos. Un artículo defectuoso se designa como un éxito. El número de éxitos es una variable aleatoria X que toma valores integrales de cero a 3. Los ocho resultados posibles y los valores correspondientes de X son

Resultado	NNN	NDN	NND	DNN	NDD	DND	DDN	DDD
x	0	1	1	1	2	2	2	3

Como los artículos se seleccionan de forma independiente y se asume que el proceso produce 25% de artículos defectuosos,

$$P(NDN) = P(N)P(D)P(N) = \left(\frac{3}{4}\right) \left(\frac{1}{4}\right) \left(\frac{3}{4}\right) = \frac{9}{64}.$$

Cálculos similares dan las probabilidades para los otros resultados posibles. La distribución de probabilidad de X es, por lo tanto,

x	0	1	2	3
$f(x)$	$\frac{27}{64}$	$\frac{27}{64}$	$\frac{9}{64}$	$\frac{1}{64}$

Distribución binomial

El número X de éxitos en n experimentos de Bernoulli se denomina **variable aleatoria binomial**. La distribución de probabilidad de esta variable aleatoria discreta se llama **distribución binomial** y sus valores se denotarán como $b(x; n, p)$, ya que dependen del número de ensayos y de la probabilidad de éxito en un ensayo dado. Por consiguiente, para la distribución de probabilidad de X el número de productos defectuosos es

$$P(X = 2) = f(2) = b\left(2; 3, \frac{1}{4}\right) = \frac{9}{64}.$$

Generalicemos ahora la ilustración anterior con el fin de obtener una fórmula para $b(x; n, p)$. Esto significa que deseamos encontrar una fórmula que dé la probabilidad de x éxitos en n ensayos para un experimento binomial. Empezamos por considerar la probabilidad de x éxitos y $n-x$ fracasos en un orden específico. Como los ensayos son independientes, podemos multiplicar todas las probabilidades que corresponden a los diferentes resultados. Cada éxito ocurre con probabilidad p y cada fracaso con probabilidad $q = 1 - p$. Por lo tanto, la probabilidad para el orden específico es $p^x q^{n-x}$. Ahora debemos determinar el número total de puntos muestrales en el experimento que tienen x éxitos y $n-x$ fracasos. Este número es igual al número de particiones de n resultados en dos grupos con x en un grupo y $n-x$ en el otro, y se escribe $\binom{n}{x}$ como se presentó en la sección 2.3. Como estas particiones son mutuamente excluyentes, sumamos las probabilidades de todas las diferentes particiones para obtener la fórmula general o simplemente multiplicamos $p^x q^{n-x}$ por $\binom{n}{x}$.

Distribución binomial Un experimento de Bernoulli puede tener como resultado un éxito con probabilidad p y un fracaso con probabilidad $q = 1 - p$. Entonces, la distribución de probabilidad de la variable aleatoria binomial X , el número de éxitos en n ensayos independientes, es

$$b(x; n, p) = \binom{n}{x} p^x q^{n-x}, \quad x = 0, 1, 2, \dots, n.$$

Observe que cuando $n = 3$ y $p = 1/4$, la distribución de probabilidad de X , el número de artículos defectuosos, se escribe como

$$b\left(x; 3, \frac{1}{4}\right) = \binom{3}{x} \left(\frac{1}{4}\right)^x \left(\frac{3}{4}\right)^{3-x}, \quad x = 0, 1, 2, 3.$$

en vez de la forma tabular de la página 144.

Ejemplo 5.1: La probabilidad de que cierta clase de componente sobreviva a una prueba de choque es de $3/4$. Calcule la probabilidad de que sobrevivan exactamente 2 de los siguientes 4 componentes que se prueben.

Solución: Si suponemos que las pruebas son independientes y $p = 3/4$ para cada una de las 4 pruebas, obtenemos

$$b\left(2; 4, \frac{3}{4}\right) = \binom{4}{2} \left(\frac{3}{4}\right)^2 \left(\frac{1}{4}\right)^2 = \binom{4!}{2! 2!} \left(\frac{3^2}{4^4}\right) = \frac{27}{128}.$$

¿De dónde proviene el nombre binomial?

La distribución binomial deriva su nombre del hecho de que los $n + 1$ términos en la expansión binomial de $(q + p)^n$ corresponden a los diversos valores de $b(x; n, p)$ para $x = 0, 1, 2, \dots, n$. Es decir,

$$\begin{aligned} (q + p)^n &= \binom{n}{0} q^n + \binom{n}{1} p q^{n-1} + \binom{n}{2} p^2 q^{n-2} + \dots + \binom{n}{n} p^n \\ &= b(0; n, p) + b(1; n, p) + b(2; n, p) + \dots + b(n; n, p). \end{aligned}$$

Dado que $p + q = 1$, vemos que

$$\sum_{x=0}^n b(x; n, p) = 1,$$

una condición que se debe cumplir para cualquier distribución de probabilidad.

Con frecuencia nos interesamos en problemas donde se necesita obtener $P(X < r)$ o $P(a \leq X \leq b)$. Las sumatorias binomiales

$$B(r, n, p) = \sum_{x=0}^r b(x; n, p)$$

se presentan en la tabla A.1 del apéndice para $n = 1, 2, \dots, 20$, para valores seleccionados de p entre 0.1 y 0.9. Ilustramos el uso de la tabla A.1 con el siguiente ejemplo.

Ejemplo 5.2: La probabilidad de que un paciente se recupere de una rara enfermedad sanguínea es de 0.4. Si se sabe que 15 personas contrajeron la enfermedad, ¿cuál es la probabilidad de que a) sobrevivan al menos 10, b) sobrevivan de 3 a 8, y c) sobrevivan exactamente 5?

Solución: Sea X el número de personas que sobreviven.

$$\begin{aligned} a) \quad P(X \geq 10) &= 1 - P(X < 10) = 1 - \sum_{x=0}^9 b(x; 15, 0.4) = 1 - 0.9662 \\ &= 0.0338 \end{aligned}$$

$$\begin{aligned} b) \quad P(3 \leq X \leq 8) &= \sum_{x=3}^8 b(x; 15, 0.4) = \sum_{x=0}^8 b(x; 15, 0.4) - \sum_{x=0}^2 b(x; 15, 0.4) \\ &= 0.9050 - 0.0271 = 0.8779 \end{aligned}$$

$$\begin{aligned} c) \quad P(X = 5) &= b(5; 15, 0.4) = \sum_{x=0}^5 b(x; 15, 0.4) - \sum_{x=0}^4 b(x; 15, 0.4) \\ &= 0.4032 - 0.2173 = 0.1859 \end{aligned}$$

Ejemplo 5.3: Una cadena grande de tiendas al detalle le compra cierto tipo de dispositivo electrónico a un fabricante, el cual le indica que la tasa de dispositivos defectuosos es de 3%.

- a) El inspector de la cadena elige 20 artículos al azar de un cargamento. ¿Cuál es la probabilidad de que haya al menos un artículo defectuoso entre estos 20?
- b) Suponga que el detallista recibe 10 cargamentos en un mes y que el inspector prueba aleatoriamente 20 dispositivos por cargamento. ¿Cuál es la probabilidad de que haya exactamente tres cargamentos que contengan al menos un dispositivo defectuoso de entre los 20 seleccionados y probados?

Solución: a) Denote con X el número de dispositivos defectuosos de los 20. Entonces X sigue una distribución $b(x; 20, 0.03)$. Por consiguiente,

$$\begin{aligned} P(X \geq 1) &= 1 - P(X = 0) = 1 - b(0; 20, 0.03) \\ &= 1 - (0.03)^0 (1 - 0.03)^{20-0} = 0.4562. \end{aligned}$$

- b) En este caso cada cargamento puede o no contener al menos un artículo defectuoso. Por lo tanto, el hecho de probar el resultado de cada cargamento puede considerarse como un experimento de Bernoulli con $p = 0.4562$ del inciso a). Si suponemos la independencia de un cargamento a otro, y si se denotamos con Y el número de cargamentos que contienen al menos un artículo defectuoso, Y sigue otra distribución bi-

nomial $b(y; 10, 0.4562)$. Por lo tanto,

$$P(Y = 3) = \binom{10}{3} 0.4562^3 (1 - 0.4562)^7 = 0.1602.$$

Áreas de aplicación

A partir de los ejemplos 5.1 a 5.3 debería quedar claro que la distribución binomial tiene aplicaciones en muchos campos científicos. Un ingeniero industrial está muy interesado en “la proporción de artículos defectuosos” en cierto proceso industrial. A menudo las medidas de control de calidad y los esquemas de muestreo para procesos se basan en la distribución binomial, la cual se aplica en cualquier situación industrial donde el resultado de un proceso es dicotómico y los resultados del proceso son independientes, y además la probabilidad de éxito se mantiene constante de una prueba a otra. La distribución binomial también se utiliza mucho en aplicaciones médicas y militares. En ambos casos un resultado de éxito o de fracaso es importante. Por ejemplo, la importancia del trabajo farmacéutico radica en poder determinar si un determinado fármaco “cura” o “no cura”; mientras que si se está probando la eficacia al lanzar un proyectil el resultado se interpretaría como “dar en el blanco” o “fallar”.

Como la distribución de probabilidad de cualquier variable aleatoria binomial depende sólo de los valores que toman los parámetros n , p y q , parecería razonable suponer que la media y la varianza de una variable aleatoria binomial también dependen de los valores que toman tales parámetros. En realidad esto es cierto, y en la demostración del teorema 5.1 derivamos fórmulas generales que se pueden utilizar para calcular la media y la varianza de cualquier variable aleatoria binomial como funciones de n , p y q .

Teorema 5.1: La media y la varianza de la distribución binomial $b(x; n, p)$ son

$$\mu = np \text{ y } \sigma^2 = npq.$$

Prueba: Representemos el resultado de la j -ésima prueba mediante una variable aleatoria de Bernoulli I_j , que toma los valores 0 y 1 con probabilidades q y p , respectivamente. Por lo tanto, en un experimento binomial el número de éxitos se escribe como la suma de las n variables indicadoras independientes. De aquí,

$$X = I_1 + I_2 + \cdots + I_n.$$

La media de cualquier I_j es $E(I_j) = (0)(q) + (1)(p) = p$. Por lo tanto, usando el corolario 4.4 de la página 131, la media de la distribución binomial es

$$\mu = E(X) = E(I_1) + E(I_2) + \cdots + E(I_n) = \underbrace{p + p + \cdots + p}_{n \text{ términos}} = np.$$

La varianza de cualquier I_j es $\sigma_{I_j}^2 = E(I_j^2) - p^2 = (0)^2(q) + (1)^2(p) - p^2 = p(1-p) = pq$. Al ampliar el corolario 4.11 al caso de n variables de Bernoulli independientes, la varianza de la distribución binomial resulta como

$$\sigma_X^2 = \sigma_{I_1}^2 + \sigma_{I_2}^2 + \cdots + \sigma_{I_n}^2 = \underbrace{pq + pq + \cdots + pq}_{n \text{ términos}} = npq.$$

Ejemplo 5.4: Se conjetura que hay impurezas en 30% del total de pozos de agua potable de cierta comunidad rural. Para obtener información sobre la verdadera magnitud del problema se determina que debe realizarse algún tipo de prueba. Como es muy costoso probar todos los pozos del área, se eligen 10 al azar para someterlos a la prueba.

- a) Si se utiliza la distribución binomial, ¿cuál es la probabilidad de que exactamente 3 pozos tengan impurezas, considerando que la conjetura es correcta?
 b) ¿Cuál es la probabilidad de que más de 3 pozos tengan impurezas?

Solución: a) Requerimos

$$b(3; 10, 0.3) = \sum_{x=0}^3 b(x; 10, 0.3) - \sum_{x=0}^2 b(x; 10, 0.3) = 0.6496 - 0.3828 = 0.2668.$$

b) En este caso $P(X > 3) = 1 - 0.6496 = 0.3504$. ┘

Ejemplo 5.5: Calcule la media y la varianza de la variable aleatoria binomial del ejemplo 5.2 y después utilice el teorema de Chebyshev (de la página 137) para interpretar el intervalo $\mu \pm 2\sigma$.

Solución: Como el ejemplo 5.2 fue un experimento binomial con $n = 15$ y $p = 0.4$, por el teorema 5.1 tenemos

$$\mu = (15)(0.4) = 6 \text{ y } \sigma^2 = (15)(0.4)(0.6) = 3.6.$$

Al tomar la raíz cuadrada de 3.6 encontramos que $\sigma = 1.897$. Por lo tanto, el intervalo que se requiere es $6 \pm (2)(1.897)$, o de 2.206 a 9.794. El teorema de Chebyshev establece que el número de pacientes recuperados, de un total de 15 que contrajeron la enfermedad, tiene una probabilidad de al menos $3/4$ de caer entre 2.206 y 9.794 o, como los datos son discretos, incluso entre 2 y 10. ┘

Hay soluciones en las que el cálculo de las probabilidades binomiales nos permitirían hacer inferencias científicas acerca de una población después de que se recaban los datos. El siguiente ejemplo es una ilustración de esto.

Ejemplo 5.6: Considere la situación del ejemplo 5.4. La idea de que el 30% de los pozos tienen impurezas es sólo una conjetura del consejo local del agua. Suponga que se eligen 10 pozos de forma aleatoria y resulta que 6 contienen impurezas. ¿Qué implica esto respecto de la conjetura? Utilice un enunciado de probabilidad.

Solución: Primero debemos preguntar: "Si la conjetura es correcta, ¿podríamos haber encontrado 6 o más pozos con impurezas?"

$$P(X \geq 6) = \sum_{x=0}^{10} b(x; 10, 0.3) - \sum_{x=0}^5 b(x; 10, 0.3) = 1 - 0.9527 = 0.0473.$$

En consecuencia, es poco probable (4.7% de probabilidad) que se encontrara que 6 o más pozos contenían impurezas si sólo 30% de ellos las contienen. Esto pone seriamente en duda la conjetura y sugiere que el problema de la impureza es mucho más grave. ┘

Como podrá darse cuenta el lector ahora, en muchas aplicaciones hay más de dos resultados posibles. Por ejemplo, en el campo de la genética el color de las crías de conejillos de Indias puede ser rojo, negro o blanco. Con frecuencia la dicotomía de "defectuoso" y "sin defectos" en casos de ingeniería es en realidad un simplificación excesiva. De hecho, a menudo hay más de dos categorías que caracterizan los artículos o las partes que salen de una línea de producción.

Experimentos multinomiales y la distribución multinomial

El experimento binomial se convierte en un experimento multinomial si cada prueba tiene más de dos resultados posibles. La clasificación de un producto fabricado como ligero, pesado o aceptable, y el registro de los accidentes en cierto crucero de acuerdo con el día de la semana, constituyen experimentos multinomiales. Extraer *con reemplazo* una carta de una baraja también es un experimento multinomial si los 4 palos son los resultados de interés.

En general, si un ensayo dado puede tener como consecuencia cualquiera de los k resultados posibles E_1, E_2, \dots, E_k con probabilidades p_1, p_2, \dots, p_k , la **distribución multinomial** dará la probabilidad de que E_1 ocurra x_1 veces, E_2 ocurra x_2 veces... y E_k ocurra x_k veces en n ensayos independientes, donde

$$x_1 + x_2 + \dots + x_k = n.$$

Denotaremos esta distribución de probabilidad conjunta como

$$f(x_1, x_2, \dots, x_k; p_1, p_2, \dots, p_k, n).$$

Salta a la vista que $p_1 + p_2 + \dots + p_k = 1$, pues el resultado de cada ensayo debe ser uno de los k resultados posibles.

Para derivar la fórmula general procedemos como en el caso binomial. Puesto que los ensayos son independientes, cualquier orden especificado que produzca x_1 resultados para E_1 , x_2 para E_2, \dots, x_k para E_k ocurrirá con probabilidad $p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$. El número total de ordenamientos que producen resultados similares para los n ensayos es igual al número de particiones de n artículos en k grupos con x_1 en el primer grupo, x_2 en el segundo grupo, ..., y x_k en el k -ésimo grupo. Esto se puede hacer en

$$\binom{n}{x_1, x_2, \dots, x_k} = \frac{n!}{x_1! x_2! \dots x_k!}$$

formas. Como todas las particiones son mutuamente excluyentes y tienen la misma probabilidad de ocurrir, obtenemos la distribución multinomial multiplicando la probabilidad para un orden específico por el número total de particiones.

Distribución multinomial Si un ensayo dado puede producir los k resultados E_1, E_2, \dots, E_k con probabilidades p_1, p_2, \dots, p_k , entonces la distribución de probabilidad de las variables aleatorias X_1, X_2, \dots, X_k , que representa el número de ocurrencias para E_1, E_2, \dots, E_k en n ensayos independientes, es

$$f(x_1, x_2, \dots, x_k; p_1, p_2, \dots, p_k, n) = \binom{n}{x_1, x_2, \dots, x_k} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k},$$

con

$$\sum_{i=1}^k x_i = n \text{ y } \sum_{i=1}^k p_i = 1.$$

La distribución multinomial deriva su nombre del hecho de que los términos de la expansión multinomial de $(p_1 + p_2 + \dots + p_k)^n$ corresponden a todos los posibles valores de $f(x_1, x_2, \dots, x_k; p_1, p_2, \dots, p_k, n)$.

Ejemplo 5.7: La complejidad de las llegadas y las salidas de los aviones en un aeropuerto es tal que a menudo se utiliza la simulación por computadora para modelar las condiciones “ideales”. Para un aeropuerto específico que tiene tres pistas se sabe que, en el escenario ideal, las probabilidades de que las pistas individuales sean utilizadas por un avión comercial que llega aleatoriamente son las siguientes:

Pista 1:	$p_1 = 2/9$
Pista 2:	$p_2 = 1/6$
Pista 3:	$p_3 = 11/18$

¿Cuál es la probabilidad de que 6 aviones que llegan al azar se distribuyan de la siguiente manera?

Pista 1:	2 aviones
Pista 2:	1 avión
Pista 3:	3 aviones

Solución: Si usamos la distribución multinomial, tenemos

$$f\left(2, 1, 3; \frac{2}{9}, \frac{1}{6}, \frac{11}{18}, 6\right) = \binom{6}{2, 1, 3} \left(\frac{2}{9}\right)^2 \left(\frac{1}{6}\right)^1 \left(\frac{11}{18}\right)^3$$

$$= \frac{6!}{2!1!3!} \cdot \frac{2^2}{9^2} \cdot \frac{1}{6} \cdot \frac{11^3}{18^3} = 0.1127.$$

Ejercicios

5.1 Una variable aleatoria X que toma los valores x_1, x_2, \dots, x_k se denomina variable aleatoria discreta uniforme si su función de masa de probabilidad es $f(x) = \frac{1}{k}$ para todas las variables x_1, x_2, \dots, x_k y 0 en cualquier otro caso. Calcule la media y la varianza de X .

5.2 Se entregan dos altavoces idénticos a 12 personas y se les pide que los escuchen para determinar si hay alguna diferencia entre ellos. Suponga que sus respuestas son simplemente conjeturas. Calcule la probabilidad de que tres personas afirmen haber detectado una diferencia entre los dos altavoces.

5.3 De un equipo de 10 empleados, y mediante la selección al azar de una etiqueta contenida en una caja que contiene 10 etiquetas numeradas del 1 al 10, se elige a uno para que supervise cierto proyecto. Calcule la fórmula para la distribución de probabilidad de X que represente el número en la etiqueta que se saca. ¿Cuál es la probabilidad de que el número que se extrae sea menor que 4?

5.4 En cierto distrito de la ciudad se establece que la causa de 75% de todos los robos es la necesidad de dinero para comprar drogas. Calcule la probabilidad de que entre los siguientes cinco casos de robo que se reporten en este distrito,

- exactamente 2 sean resultado de la necesidad de dinero para comprar drogas;
- a lo sumo 3 resulten de la necesidad de dinero para comprar drogas.

5.5 De acuerdo con *Chemical Engineering Progress* (noviembre de 1990), aproximadamente 30% de todas las fallas de operación en las tuberías de plantas químicas son ocasionadas por errores del operador.

- ¿Cuál es la probabilidad de que de las siguientes 20 fallas en las tuberías al menos 10 se deban a un error del operador?
- ¿Cuál es la probabilidad de que no más de 4 de 20 fallas se deban a un error del operador?
- Suponga que, para una planta específica, de la muestra aleatoria de 20 de tales fallas exactamente 5 son errores de operación. ¿Considera que la cifra de 30% anterior se aplique a esta planta? Comente su respuesta.

5.6 De acuerdo con una encuesta de la *Administrative Management Society*, la mitad de las empresas estadounidenses da a sus empleados 4 semanas de vacaciones después de 15 años de servicio en la empresa. Calcule la probabilidad de que, de 6 empresas encuestadas al azar, el número que da a sus empleados 4 semanas de vacaciones después de 15 años de servicio es

- cualquiera entre 2 y 5;
- menor que 3.

5.7 Un destacado médico afirma que el 70% de las personas con cáncer de pulmón son fumadores empedernidos. Si su aseveración es correcta,

- calcule la probabilidad de que de 10 de estos pacientes, que ingresaron recientemente a un hospital, menos de la mitad sean fumadores empedernidos;

b) calcule la probabilidad de que de 20 de estos pacientes, que ingresaron recientemente a un hospital, menos de la mitad sean fumadores empedernidos.

5.8 De acuerdo con un estudio publicado por un grupo de sociólogos de la Universidad de Massachusetts, aproximadamente 60% de los consumidores de Valium en el estado de Massachusetts empezaron a consumirlo a causa de problemas psicológicos. Calcule la probabilidad de que entre los siguientes 8 consumidores entrevistados de este estado,

- exactamente 3 comenzaron a consumir Valium por problemas psicológicos;
- al menos 5 comenzaron a consumir Valium por problemas que no fueron psicológicos.

5.9 Al probar cierta clase de neumático para camión en un terreno accidentado, se encuentra que el 25% de los camiones no completan la prueba de recorrido sin ponchaduras. De los siguientes 15 camiones probados, calcule la probabilidad de que

- de 3 a 6 tengan ponchaduras;
- menos de 4 tengan ponchaduras;
- más de 5 tengan ponchaduras.

5.10 Según un informe de la revista *Parade*, una encuesta a nivel nacional, realizada por la Universidad de Michigan con estudiantes universitarios de último año, reveló que casi 70% desaprueban el consumo diario de marihuana. Si se seleccionan 12 estudiantes de último año al azar y se les pide su opinión, calcule la probabilidad de que el número de los que desaprueban el consumo diario de marihuana sea

- cualquiera entre 7 y 9;
- 5 a lo sumo;
- no menos de 8.

5.11 La probabilidad de que un paciente se recupere de una delicada operación de corazón es 0.9. ¿Cuál es la probabilidad de que exactamente 5 de los siguientes 7 pacientes intervenidos sobrevivan?

5.12 Un ingeniero de control de tráfico reporta que 75% de los vehículos que pasan por un punto de verificación son de ese estado. ¿Cuál es la probabilidad de que menos de 4 de los siguientes 9 vehículos sean de otro estado?

5.13 Un estudio a nivel nacional que examinó las actitudes hacia los antidepresivos reveló que aproximadamente 70% de los encuestados cree que "los antidepresivos en realidad no curan nada, sólo disfrazan el problema real". De acuerdo con este estudio, ¿cuál es la probabilidad de que al menos 3 de las siguientes 5 personas seleccionadas al azar tengan esta opinión?

5.14 El porcentaje de victorias que consiguió el equipo de baloncesto los Toros de Chicago para pasar a las

finales en la temporada 1996-97 fue de 87.7. Redondee 87.7 a 90 para poder utilizar la tabla A.1.

- ¿Cuál es la probabilidad de que los Toros logren una victoria aplastante (4-0) en la serie final de 7 juegos?
- ¿Cuál es la probabilidad de que los Toros ganen la serie inicial?
- ¿Qué suposición importante se hace al responder los incisos a) y b)?

5.15 Se sabe que 60% de los ratones inoculados con un suero quedan protegidos contra cierta enfermedad. Si se inoculan 5 ratones, calcule la probabilidad de que

- ninguno contraiga la enfermedad;
- menos de 2 contraigan la enfermedad;
- más de 3 contraigan la enfermedad.

5.16 Suponga que los motores de un avión operan de forma independiente y que tienen una probabilidad de falla de 0.4. Se supone que un avión tiene un vuelo seguro si funcionan al menos la mitad de sus motores. Si un avión tiene 4 motores y otro tiene 2, ¿cuál de los dos tiene la probabilidad más alta de un vuelo exitoso?

5.17 Si X representa el número de personas del ejercicio 5.13 que creen que los antidepresivos no curan sino que sólo disfrazan el problema real, calcule la media y la varianza de X si se seleccionan al azar 5 personas.

5.18 a) ¿Cuántos de los 15 camiones del ejercicio 5.9 esperarían que tuvieran ponchaduras?

- ¿Cuál es la varianza del número de ponchaduras de los 15 camiones? ¿Qué significado tiene eso?

5.19 Un estudiante que conduce hacia su escuela encuentra un semáforo, el cual permanece verde por 35 segundos, amarillo cinco segundos y rojo 60 segundos. Suponga que toda la semana el estudiante recorre el camino a la escuela entre las 8:00 y las 8:30 a.m. Sea X_1 el número de veces que encuentra una luz verde, X_2 el número de veces que encuentra una luz amarilla y X_3 el número de veces que encuentra una luz roja. Calcule la distribución conjunta de X_1 , X_2 y X_3 .

5.20 Según el diario *USA Today* (18 de marzo de 1997), de 4 millones de integrantes de la fuerza laboral, 5.8% resultó positivo en una prueba de drogas. De los que dieron positivo, 22.5% consumían cocaína y 54.4% consumían marihuana.

- ¿Cuál es la probabilidad de que de 10 trabajadores que dieron positivo, 2 sean usuarios de cocaína, 5 de marihuana y 3 de otras drogas?
- ¿Cuál es la probabilidad de que de 10 trabajadores que dieron positivo, todos sean consumidores de marihuana?
- ¿Cuál es la probabilidad de que de 10 trabajadores que dieron positivo, ninguno consuma cocaína?

5.21 La superficie de un tablero circular para dardos tiene un pequeño círculo central llamado diana y 20 regiones en forma de rebanada de pastel numeradas del 1 al 20. Asimismo, cada una de estas regiones está dividida en tres partes, de manera que una persona que lanza un dardo que cae en un número específico obtiene una puntuación igual al valor del número, el doble del número o el triple de éste, dependiendo de en cuál de las tres partes caiga el dardo. Si una persona tiene una probabilidad de 0.01 de acertar a la diana, una probabilidad de 0.10 de acertar un doble, una probabilidad de 0.05 de acertar un triple y una probabilidad de 0.02 de no acertar al tablero, ¿cuál es la probabilidad de que 7 lanzamientos den como resultado ninguna diana, ningún triple, dos dobles y una vez fuera del tablero?

5.22 De acuerdo con la teoría genética, cierta cruce de conejillos de Indias tendrá crías rojas, negras y blancas en la proporción 8:4:4. Calcule la probabilidad de que de 8 crías, 5 sean rojas, 2 negras y 1 blanca.

5.23 Las probabilidades de que un delegado llegue a cierta convención en avión, autobús, automóvil o tren son de 0.4, 0.2, 0.3 y 0.1, respectivamente. ¿Cuál es la probabilidad de que, de 9 delegados que asisten a esta convención seleccionados al azar, 3 lleguen en avión, 3 en autobús, 1 en automóvil y 2 en tren?

5.24 Un ingeniero de seguridad afirma que sólo 40% de los trabajadores utilizan cascos de seguridad cuando comen en el lugar de trabajo. Suponga que esta afirmación es cierta y calcule la probabilidad de que 4 de 6 trabajadores elegidos al azar utilicen sus cascos mientras comen en el lugar de trabajo.

5.25 Suponga que para un embarque muy grande de circuitos integrados, la probabilidad de que falle cualquiera de ellos es de 0.10. Suponga que se cumplen los supuestos en que se basan las distribuciones binomiales y calcule la probabilidad de que en una muestra aleatoria de 20 fallen, a lo sumo, 3 chips integrados.

5.26 Suponga que 6 de 10 accidentes automovilísticos se deben principalmente a que no se respeta el límite de velocidad y calcule la probabilidad de que, de 8 accidentes automovilísticos, 6 se deban principalmente a una violación del límite de velocidad

- mediante el uso de la fórmula para la distribución binomial;
- usando la tabla A.1.

5.27 Si una bombilla fluorescente tiene una probabilidad de 0.9 de tener una vida útil de al menos 800 horas, calcule las probabilidades de que, de 20 bombillas fluorescentes,

- exactamente 18 tengan una vida útil de al menos 800 horas;
- al menos 15 tengan una vida útil de al menos 800 horas;
- al menos 2 *no* tengan una vida útil de al menos 800 horas.

5.28 Un fabricante sabe que, en promedio, 20% de los tostadores eléctricos producidos requerirá reparaciones durante el primer año posterior a su venta. Suponga que se seleccionan al azar 20 tostadores y calcule los números x y y adecuados tales que

- la probabilidad de que al menos x de ellos requieran reparaciones sea menor que 0.5;
- la probabilidad de que al menos y de ellos *no* requieran reparaciones sea mayor que 0.8.

5.3 Distribución hipergeométrica

La manera más simple de ver la diferencia entre la distribución binomial de la sección 5.2 y la distribución hipergeométrica consiste en observar la forma en que se realiza el muestreo. Los tipos de aplicaciones de la distribución hipergeométrica son muy similares a los de la distribución binomial. Nos interesa el cálculo de probabilidades para el número de observaciones que caen en una categoría específica. Sin embargo, la distribución binomial requiere que los ensayos sean independientes. Por consiguiente, si se aplica esta distribución, digamos, tomando muestras de un lote de artículos (barajas, lotes de artículos producidos), el muestreo se debe efectuar **reemplazando** cada artículo después de observarlo. Por otro lado, la distribución hipergeométrica *no* requiere independencia y se basa en el muestreo que se realiza **sin reemplazo**.

Las aplicaciones de la distribución hipergeométrica se encuentran en muchos campos, sobre todo en el muestreo de aceptación, las pruebas electrónicas y los controles de calidad. Evidentemente, en muchos de estos campos el muestreo se realiza a expensas del artículo que se prueba; es decir, el artículo se destruye, por lo que no se puede

reemplazar en la muestra. Por consiguiente, el muestreo sin reemplazo es necesario. Utilizaremos un caso simple con barajas para nuestro primer ejemplo.

Si deseamos calcular la probabilidad de obtener 3 cartas rojas en 5 extracciones de una baraja ordinaria de 52 cartas, la distribución binomial de la sección 5.2 no se aplica a menos que cada carta se reemplace y que el paquete se revuelva antes de extraer la siguiente carta. Para resolver el problema del muestreo sin reemplazo volvamos a plantear el problema. Si se sacan 5 cartas al azar, nos interesa la probabilidad de seleccionar 3 cartas rojas de las 26 disponibles y 2 de las 26 cartas negras de que dispone la baraja. Hay $\binom{26}{3}$ formas de seleccionar 3 cartas rojas, y para cada una de estas formas podemos elegir 2 cartas negras de $\binom{26}{2}$ maneras. Por lo tanto, el número total de formas de seleccionar 3 cartas rojas y 2 negras en 5 extracciones es el producto $\binom{26}{3}\binom{26}{2}$. El número total de formas de seleccionar cualesquiera 5 cartas de las 52 disponibles es $\binom{52}{5}$. En consecuencia, la probabilidad de seleccionar 5 cartas sin reemplazo, de las cuales 3 sean rojas y 2 negras está dada por

$$\frac{\binom{26}{3}\binom{26}{2}}{\binom{52}{5}} = \frac{(26!/3!23!)(26!/2!24!)}{52!/5!47!} = 0.3251.$$

En general, nos interesa la probabilidad de seleccionar x éxitos de los k artículos considerados éxitos y $n - x$ fracasos de los $N - k$ artículos que se consideran fracasos cuando una muestra aleatoria de tamaño n se selecciona de N artículos. Esto se conoce como un **experimento hipergeométrico**; es decir, aquel que posee las siguientes dos propiedades:

1. De un lote de N artículos se selecciona una muestra aleatoria de tamaño n sin reemplazo.
2. k de los N artículos se pueden clasificar como éxitos y $N - k$ se clasifican como fracasos.

El número X de éxitos de un experimento hipergeométrico se denomina **variable aleatoria hipergeométrica**. En consecuencia, la distribución de probabilidad de la variable hipergeométrica se conoce como **distribución hipergeométrica**, y sus valores se denotan con $h(x; N, n, k)$, ya que dependen del número de éxitos k en el conjunto N del que seleccionamos n artículos.

Distribución hipergeométrica en el muestreo de aceptación

Como en el caso de la distribución binomial, la distribución hipergeométrica se aplica en el muestreo de aceptación, donde se toman muestras del material o las partes de los lotes con el fin de determinar si se acepta o no el lote completo.

Ejemplo 5.8: Una parte específica que se utiliza como dispositivo de inyección se vende en lotes de 10. El productor considera que el lote es aceptable si no tiene más de un artículo defectuoso. Un plan de muestreo incluye un muestreo aleatorio y la prueba de 3 de cada 10 partes. Si ninguna de las 3 está defectuosa, se acepta el lote. Comente acerca de la utilidad de este plan.

Solución: Supongamos que el lote es verdaderamente **inaceptable** (es decir, que 2 de cada 10 partes están defectuosas). La probabilidad de que el plan de muestreo considere que el lote aceptable es

$$P(X = 0) = \frac{\binom{2}{0}\binom{8}{3}}{\binom{10}{3}} = 0.467.$$

Por consiguiente, si el lote es realmente inaceptable porque 2 partes están defectuosas, este plan de muestreo permitirá que se acepte aproximadamente 47% de las veces. Como resultado, este plan debería considerarse inadecuado. \blacksquare

Hagamos una generalización para calcular una fórmula para $h(x; N, n, k)$. El número total de muestras de tamaño n elegidas de N artículos es $\binom{N}{n}$. Se supone que estas muestras tienen la misma probabilidad. Hay $\binom{k}{x}$ formas de seleccionar x éxitos de los k disponibles, y por cada una de estas formas podemos elegir $n-x$ fracasos en formas $\binom{N-k}{n-x}$. De esta manera, el número total de muestras favorables entre las $\binom{N}{n}$ muestras posibles, está dado por $\binom{k}{x} \binom{N-k}{n-x}$. En consecuencia, tenemos la siguiente definición.

Distribución hipergeométrica La distribución de probabilidad de la variable aleatoria hipergeométrica X , el número de éxitos en una muestra aleatoria de tamaño n que se selecciona de N artículos, en los que k se denomina éxito y $N-k$ fracaso, es

$$h(x; N, n, k) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}}, \quad \text{máx}\{0, n - (N - k)\} \leq x \leq \text{mín}\{n, k\}.$$

El rango de x puede determinarse mediante los tres coeficientes binomiales en la definición, donde x y $n-x$ no son más que k y $N-k$; respectivamente; y ambos no pueden ser menores que 0. Por lo general, cuando tanto k (el número de éxitos) como $N-k$ (el número de fracasos) son mayores que el tamaño de la muestra n , el rango de una variable aleatoria hipergeométrica será $x = 0, 1, \dots, n$.

Ejemplo 5.9: Lotes con 40 componentes cada uno que contengan 3 o más defectuosos se consideran inaceptables. El procedimiento para obtener muestras del lote consiste en seleccionar 5 componentes al azar y rechazar el lote si se encuentra un componente defectuoso. ¿Cuál es la probabilidad de, que en la muestra, se encuentre exactamente un componente defectuoso, si en todo el lote hay 3 defectuosos?

Solución: Si utilizamos la distribución hipergeométrica con $n = 5$, $N = 40$, $k = 3$ y $x = 1$, encontramos que la probabilidad de obtener un componente defectuoso es

$$h(1; 40, 5, 3) = \frac{\binom{3}{1} \binom{37}{4}}{\binom{40}{5}} = 0.3011.$$

De nueva cuenta este plan no es adecuado porque sólo 30% de las veces detecta un lote malo (con 3 componentes defectuosos). \blacksquare

Teorema 5.2: La media y la varianza de la distribución hipergeométrica $h(x; N, n, k)$ son

$$\mu = \frac{nk}{N} \quad \text{y} \quad \sigma^2 = \frac{N-n}{N-1} \cdot n \cdot \frac{k}{N} \left(1 - \frac{k}{N}\right).$$

La demostración para la media se muestra en el apéndice A.24.

Ejemplo 5.10: Volvamos a investigar el ejemplo 3.4 de la página 83. La finalidad de este ejemplo fue ilustrar el concepto de una variable aleatoria y el espacio muestral correspondiente. En el ejemplo tenemos un lote de 100 artículos, de los cuales 12 están defectuosos. ¿Cuál es la probabilidad de que haya 3 defectuosos en una muestra de 10?

Solución: Si utilizamos la función de probabilidad hipergeométrica, tenemos

$$h(3; 100, 10, 12) = \frac{\binom{12}{3} \binom{88}{7}}{\binom{100}{10}} = 0.08.$$

Ejemplo 5.11: Calcule la media y la varianza de la variable aleatoria del ejemplo 5.9, y después utilice el teorema de Chebyshev para interpretar el intervalo $\mu \pm 2\sigma$.

Solución: Como el ejemplo 5.9 fue un experimento hipergeométrico con $N = 40$, $n = 5$ y $k = 3$, usando el teorema 5.2, tenemos

$$\mu = \frac{(5)(3)}{40} = \frac{3}{8} = 0.375,$$

y

$$\sigma^2 = \left(\frac{40-5}{39}\right) (5) \left(\frac{3}{40}\right) \left(1 - \frac{3}{40}\right) = 0.3113.$$

Si calculamos la raíz cuadrada de 0.3113, encontramos que $\sigma = 0.558$. Por lo tanto, el intervalo que se requiere es $0.375 \pm (2)(0.558)$, o de -0.741 a 1.491 . El teorema de Chebyshev establece que el número de componentes defectuosos que se obtienen cuando, de un lote de 40 componentes, se seleccionan 5 al azar, de los cuales 3 están defectuosos, tiene una probabilidad de al menos $3/4$ de caer entre -0.741 y 1.491 . Esto es, al menos tres cuartas partes de las veces los 5 componentes incluirán menos de 2 defectuosos.

Relación con la distribución binomial

En este capítulo examinamos varias distribuciones discretas importantes que tienen diversas aplicaciones. Muchas de estas distribuciones se relacionan bien entre sí. El estudiante novato debería tener una clara comprensión de tales relaciones. Existe una relación interesante entre las distribuciones hipergeométrica y binomial. Como se esperaría, si n es pequeña comparada con N , la naturaleza de los N artículos cambia muy poco en cada prueba. Así, cuando n es pequeña en comparación con N , se puede utilizar una distribución binomial para aproximar la distribución hipergeométrica. De hecho, por regla general la aproximación es buena cuando $n/N \leq 0.05$.

Por lo tanto, la cantidad k/N desempeña el papel del parámetro binomial p y, como consecuencia, la distribución binomial se podría considerar una versión de población grande de la distribución hipergeométrica. La media y la varianza entonces se obtienen de las fórmulas

$$\mu = np = \frac{nk}{N} \quad \text{y} \quad \sigma^2 = npq = n \cdot \frac{k}{N} \left(1 - \frac{k}{N}\right).$$

Al comparar estas fórmulas con las del teorema 5.2, vemos que la media es la misma, mientras que la varianza difiere por un factor de corrección de $(N-n)/(N-1)$, que es insignificante cuando n es pequeña en relación con N .

Ejemplo 5.12: Un fabricante de neumáticos para automóvil reporta que de un cargamento de 5000 piezas que se mandan a un distribuidor local, 1000 están ligeramente manchadas. Si se compran al azar 10 de estos neumáticos al distribuidor, ¿cuál es la probabilidad de que exactamente 3 estén manchados?

Solución: Como $N = 5000$ es grande con respecto a la muestra de tamaño $n = 10$, nos aproximaremos a la probabilidad deseada usando la distribución binomial. La probabilidad de obtener un neumático manchado es 0.2. Por lo tanto, la probabilidad de obtener exactamente 3 manchados es

$$h(3; 5000, 10, 1000) \approx b(3; 10, 0.2) = 0.8791 - 0.6778 = 0.2013.$$

Por otro lado, la probabilidad exacta es $h(3; 5000, 10, 1000) = 0.2015$.

La distribución hipergeométrica se puede extender para tratar el caso donde los N artículos se pueden dividir en k celdas A_1, A_2, \dots, A_k con a_1 elementos en la primera celda, a_2 en la segunda, ..., a_k elementos en la k -ésima celda. Lo que nos interesa ahora es la probabilidad de que una muestra aleatoria de tamaño n produzca x_1 elementos de A_1 , x_2 elementos de A_2 , ..., y x_k elementos de A_k . Representemos esta probabilidad mediante

$$f(x_1, x_2, \dots, x_k; a_1, a_2, \dots, a_k, N, n).$$

Para obtener una fórmula general observamos que el número total de muestras de tamaño n que se pueden elegir a partir de N artículos es aún $\binom{N}{n}$. Hay $\binom{a_1}{x_1}$ formas de seleccionar x_1 artículos de los que hay en A_1 , y para cada uno de éstos podemos elegir x_2 de los de A_2 en $\binom{a_2}{x_2}$ formas. Por lo tanto, podemos seleccionar x_1 artículos de A_1 y x_2 artículos de A_2 en $\binom{a_1}{x_1} \binom{a_2}{x_2}$ formas. Si continuamos de esta forma, podemos seleccionar todos los n artículos que constan de x_1 de A_1 , x_2 de A_2 , ..., y x_k de A_k en

$$\binom{a_1}{x_1} \binom{a_2}{x_2} \cdots \binom{a_k}{x_k} \text{ formas.}$$

La distribución de probabilidad que se requiere se define ahora como sigue.

Distribución hipergeométrica multivariada Si N artículos se pueden dividir en las k celdas A_1, A_2, \dots, A_k con a_1, a_2, \dots, a_k elementos, respectivamente, entonces la distribución de probabilidad de las variables aleatorias X_1, X_2, \dots, X_k , que representan el número de elementos que se seleccionan de A_1, A_2, \dots, A_k en una muestra aleatoria de tamaño n , es

$$f(x_1, x_2, \dots, x_k; a_1, a_2, \dots, a_k, N, n) = \frac{\binom{a_1}{x_1} \binom{a_2}{x_2} \cdots \binom{a_k}{x_k}}{\binom{N}{n}},$$

$$\text{con } \sum_{i=1}^k x_i = n \text{ y } \sum_{i=1}^k a_i = N.$$

Ejemplo 5.13: Se usa un grupo de 10 individuos para un estudio de caso biológico. El grupo contiene 3 personas con sangre tipo O, 4 con sangre tipo A y 3 con tipo B. ¿Cuál es la probabilidad de que una muestra aleatoria de 5 contenga 1 persona con sangre tipo O, 2 personas con tipo A y 2 personas con tipo B?

Solución: Si se utiliza la extensión de la distribución hipergeométrica con $x_1 = 1, x_2 = 2, x_3 = 2, a_1 = 3, a_2 = 4, a_3 = 3, N = 10$ y $n = 5$, vemos que la probabilidad que se desea es

$$f(1, 2, 2; 3, 4, 3, 10, 5) = \frac{\binom{3}{1} \binom{4}{2} \binom{3}{2}}{\binom{10}{5}} = \frac{3}{14}.$$

Ejercicios

5.29 El dueño de una casa planta 6 bulbos seleccionados al azar de una caja que contiene 5 bulbos de tulipán y 4 de narciso. ¿Cuál es la probabilidad de que plante 2 bulbos de narciso y 4 de tulipán?

5.30 Para evitar la detección en la aduana, un viajero coloca 6 comprimidos con narcóticos en una botella que contiene 9 píldoras de vitamina que aparentemente son similares. Si el oficial de la aduana selecciona 3 de las tabletas al azar para su análisis, ¿cuál es la probabilidad de que el viajero sea arrestado por posesión ilegal de narcóticos?

5.31 Se selecciona al azar un comité de 3 personas a partir de 4 médicos y 2 enfermeras. Escriba una fórmula para la distribución de probabilidad de la variable aleatoria X que representa el número de médicos en el comité. Calcule $P(2 \leq X \leq 3)$.

5.32 De un lote de 10 misiles, se seleccionan 4 al azar y se disparan. Si el lote contiene 3 misiles defectuosos que no pueden dispararse, ¿cuál es la probabilidad de que

- los 4 puedan dispararse?
- a lo sumo fallen 2?

5.33 Si de una baraja ordinaria de 52 cartas, se toman 7 y se reparten, ¿cuál es la probabilidad de que

- exactamente 2 de ellas sean cartas de figuras?
- al menos 1 de ellas sea una reina?

5.34 ¿Cuál es la probabilidad de que una camarera se rehúse a servir bebidas alcohólicas a sólo 2 menores si verifica al azar 5 identificaciones de 9 estudiantes, de los cuales 4 son menores de edad?

5.35 Una empresa está interesada en evaluar su procedimiento de inspección actual para embarques de 50 artículos idénticos. El procedimiento consiste en tomar una muestra de 5 artículos y aceptar el embarque si no se encuentran más de 2 defectuosos. ¿Qué proporción de embarques con 20% de artículos defectuosos se aceptará?

5.36 Una empresa de manufactura utiliza un esquema de aceptación para los artículos de una línea de producción antes de que se embarquen. El plan tiene dos etapas. Se preparan cajas de 25 artículos para su embarque y se prueba una muestra de 3 en busca de defectuosos. Si se encuentra alguno defectuoso, se regresa toda la caja para verificar el 100% de ellos. Si no se encuentran artículos defectuosos, la caja se embarca.

- ¿Cuál es la probabilidad de que se embarque una caja que contiene 3 defectuosos?
- ¿Cuál es la probabilidad de que se regrese para su revisión una caja que contenga sólo un artículo defectuoso?

5.37 Suponga que la empresa fabricante del ejercicio 5.36 decide cambiar su esquema de aceptación. Con el nuevo esquema un inspector toma un artículo al azar, lo inspecciona y después lo regresa a la caja; un segundo inspector hace lo mismo. Finalmente, un tercer inspector lleva a cabo el mismo procedimiento. Si cualquiera de los tres encuentra un artículo defectuoso, la caja no se embarca. Responda los incisos del ejercicio 5.36 con este nuevo plan.

5.38 De los 150 empleados de hacienda en una ciudad grande, sólo 30 son mujeres. Suponga que se eligen al azar 10 de los empleados para que proporcionen asesoría gratuita sobre declaraciones de impuestos a los residentes de esta ciudad; utilice la aproximación binomial a la distribución hipergeométrica para calcular la probabilidad de que se seleccionen al menos 3 mujeres.

5.39 Una ciudad vecina considera entablar una demanda de anexión en contra de una subdivisión del condado de 1200 residencias. Si los ocupantes de la mitad de las residencias objetan la anexión, ¿cuál es la probabilidad de que en una muestra aleatoria de 10 residencias al menos 3 estén a favor de la anexión?

5.40 Se estima que 4000 de los 10,000 residentes con derecho al voto de una ciudad están en contra de un nuevo impuesto sobre las ventas. Si se seleccionan al azar 15 votantes y se les pide su opinión, ¿cuál es la probabilidad de que a lo sumo 7 estén a favor del nuevo impuesto?

5.41 Una encuesta a nivel nacional, realizada por la Universidad de Michigan a 17,000 estudiantes universitarios de último año, revela que casi 70% desapruueba el consumo diario de marihuana. Si se seleccionan al azar 18 de tales estudiantes y se les pide su opinión, ¿cuál es la probabilidad de que más de 9 pero menos de 14 desaprueben el consumo de marihuana?

5.42 Calcule la probabilidad de que si le toca una mano de bridge de 13 cartas, ésta incluya 5 espadas, 2 corazones, 3 diamantes y 3 tréboles.

5.43 Un club de estudiantes extranjeros tiene como miembros a 2 canadienses, 3 japoneses, 5 italianos y 2 alemanes. Si se selecciona al azar un comité de 4, calcule la probabilidad de que

- todas las nacionalidades estén representadas;
- todas las nacionalidades estén representadas, excepto la italiana.

5.44 Una urna contiene 3 bolas verdes, 2 azules y 4 rojas. Calcule la probabilidad de que, en una muestra aleatoria de 5 bolas, se seleccionen las 2 bolas azules y al menos una roja.

5.45 A menudo los biólogos que estudian un ambiente específico etiquetan y liberan a sujetos con el fin de estimar el tamaño de la población o la prevalencia de ciertas características en ella. Los biólogos capturan a 10 animales de una especie que se piensa extinta (o casi extinta), los etiquetan y los liberan en cierta región. Después de un periodo seleccionan en la región una muestra aleatoria de 15 animales de ese tipo. ¿Cuál es la probabilidad de que 5 de los animales seleccionados estén etiquetados, si hay 25 animales de este tipo en la región?

5.46 Una empresa grande tiene un sistema de inspección para los lotes de compresores pequeños que compra a los vendedores. Un lote típico contiene 15 compresores. En el sistema de inspección se selecciona una muestra aleatoria de 5 compresores para someterlos a prueba. Suponga que en el lote de 15 hay 2 defectuosos.

- ¿Cuál es la probabilidad de que en una muestra determinada haya un compresor defectuoso?
- ¿Cuál es la probabilidad de que la inspección descubra los 2 compresores defectuosos?

5.47 Una fuerza de tareas gubernamental sospecha que algunas fábricas infringen los reglamentos federales contra la contaminación ambiental en lo que se refiere a la descarga de cierto tipo de producto. Veinte empresas están bajo sospecha pero no todas se pueden inspeccionar. Suponga que 3 de las empresas infringen los reglamentos.

- ¿Cuál es la probabilidad de que si se inspeccionan 5 empresas no se encuentre ninguna infracción?
- ¿Cuál es la probabilidad de que la inspección de 5 empresas descubra a 2 que infringen el reglamento?

5.48 Una máquina llena 10,000 latas de bebida gaseosa por hora, de entre las cuales 300 resultan con el líquido incompleto. Cada hora se elige al azar una muestra de 30 latas y se verifica el número de onzas de gaseosa que contiene cada una. Denote con X el número de latas seleccionadas con llenado insuficiente. Encuentre la probabilidad de encontrar al menos una de las latas muestreadas con llenado insuficiente.

5.4 Distribuciones binomial negativa y geométrica

Consideremos un experimento con las mismas propiedades de un experimento binomial sólo que en este caso las pruebas se repetirán hasta que ocurra un número *fijo* de éxitos. Por lo tanto, en vez de encontrar la probabilidad de x éxitos en n pruebas, donde n es fija, ahora nos interesa la probabilidad de que ocurra el k -ésimo éxito en la x -ésima prueba. Los experimentos de este tipo se llaman **experimentos binomiales negativos**.

Como ejemplo, considere el uso de un medicamento que se sabe que es eficaz en el 60% de los casos en que se utiliza. El uso del medicamento se considerará un éxito si proporciona algún grado de alivio al paciente. Nos interesa calcular la probabilidad de que el quinto paciente que experimente alivio sea el séptimo paciente en recibir el medicamento en una semana determinada. Si designamos un éxito con E y un fracaso con F , un orden posible para alcanzar el resultado que se desea es $EFEEFE$, que ocurre con la siguiente probabilidad

$$(0.6)(0.4)(0.6)(0.6)(0.6)(0.4)(0.6) = (0.6)^5(0.4)^2.$$

Podríamos listar todos los posibles ordenamientos reacomodando las F y las E , con excepción del último resultado, que debe ser el quinto éxito. El número total de ordenamientos posibles es igual al número de particiones de los primeros 6 ensayos en 2 grupos con dos fracasos asignados a un grupo y 4 éxitos asignados al otro grupo. Esto se puede realizar en $\binom{6}{4} = 15$ formas mutuamente excluyentes. Por lo tanto, si X representa el resultado en el que ocurre el quinto éxito, entonces

$$P(X = 7) = \binom{6}{4}(0.6)^5(0.4)^2 = 0.1866.$$

¿Cuál es la variable aleatoria binomial negativa?

El número X de ensayos necesarios para generar k éxitos en un experimento binomial negativo se denomina **variable aleatoria binomial negativa** y su distribución de probabi-

lidad se llama **distribución binomial negativa**. Dado que sus probabilidades dependen del número de éxitos deseados y de la probabilidad de un éxito en un ensayo dado, denotaremos ambas probabilidades con el símbolo $b^*(x; k, p)$. Para obtener la fórmula general para $b^*(x; k, p)$, considere la probabilidad de un éxito en el x -ésimo ensayo precedido por $k-1$ éxitos y $x-k$ fracasos en un orden específico. Como los ensayos son independientes podemos multiplicar todas las probabilidades que corresponden a cada resultado deseado. La probabilidad de que ocurra un éxito es p y la probabilidad de que ocurra un fracaso es $q = 1 - p$. Por lo tanto, la probabilidad para el orden específico, que termina en un éxito, es

$$p^{k-1} q^{x-k} p = p^k q^{x-k}.$$

El número total de puntos muestrales en el experimento que termina en un éxito, después de la ocurrencia de $k-1$ éxitos y $x-k$ fracasos en cualquier orden, es igual al número de particiones de $x-1$ ensayos en dos grupos con $k-1$ éxitos, que corresponden a un grupo, y $x-k$ fracasos, que corresponden al otro grupo. Este número se especifica con el término $\binom{x-1}{k-1}$, cada uno es mutuamente excluyente y tiene las mismas probabilidades de ocurrir $p^k q^{x-k}$. Obtenemos la fórmula general multiplicando $p^k q^{x-k}$ por $\binom{x-1}{k-1}$.

Distribución binomial negativa Si ensayos independientes repetidos pueden dar como resultado un éxito con probabilidad p y un fracaso con probabilidad $q = 1 - p$, entonces la distribución de probabilidad de la variable aleatoria X , el número del ensayo en el que ocurre el k -ésimo éxito, es

$$b^*(x; k, p) = \binom{x-1}{k-1} p^k q^{x-k}, \quad x = k, k+1, k+2, \dots$$

Ejemplo 5.14: En la serie de campeonato de la NBA (National Basketball Association), el equipo que gane 4 de 7 juegos será el ganador. Suponga que los equipos A y B se enfrentan en los juegos de campeonato y que el equipo A tiene una probabilidad de 0.55 de ganarle al equipo B .

- ¿Cuál es la probabilidad de que el equipo A gane la serie en 6 juegos?
- ¿Cuál es la probabilidad de que el equipo A gane la serie?
- Si ambos equipos se enfrentaran en la eliminatoria de una serie regional y el triunfador fuera el que ganara 3 de 5 juegos, ¿cuál es la probabilidad de que el equipo A gane la serie?

Solución: a) $b^*(6; 4, 0.55) = \binom{5}{3} 0.55^4 (1 - 0.55)^{6-4} = 0.1853$.

b) P (el equipo A gana la serie de campeonato) es

$$\begin{aligned} b^*(4; 4, 0.55) + b^*(5; 4, 0.55) + b^*(6; 4, 0.55) + b^*(7; 4, 0.55) \\ = 0.0915 + 0.1647 + 0.1853 + 0.1668 = 0.6083. \end{aligned}$$

c) P (el equipo A gana la eliminatoria) es

$$\begin{aligned} b^*(3; 3, 0.55) + b^*(4; 3, 0.55) + b^*(5; 3, 0.55) \\ = 0.1664 + 0.2246 + 0.2021 = 0.5931. \end{aligned}$$

La distribución binomial negativa deriva su nombre del hecho de que cada término de la expansión de $p^k(1-q)^{-k}$ corresponde a los valores de $b^*(x; k, p)$ para $x = k, k+1, k+2, \dots$. Si consideramos el caso especial de la distribución binomial negativa, donde $k = 1$, tenemos una distribución de probabilidad para el número de ensayos que se requieren para un solo éxito. Un ejemplo sería lanzar una moneda hasta que salga una cara. Nos podemos interesar en la probabilidad de que la primera cara resulte en el cuarto lanzamiento. En este caso la distribución binomial negativa se reduce a la forma

$$b^*(x; 1, p) = pq^{x-1}, \quad x = 1, 2, 3, \dots$$

Como los términos sucesivos constituyen una progresión geométrica, se acostumbra referirse a este caso especial como **distribución geométrica** y denotar sus valores con $g(x; p)$.

Distribución geométrica Si pruebas independientes repetidas pueden tener como resultado un éxito con probabilidad p y un fracaso con probabilidad $q = 1 - p$, entonces la distribución de probabilidad de la variable aleatoria X , el número de la prueba en el que ocurre el primer éxito, es

$$g(x; p) = pq^{x-1}, \quad x = 1, 2, 3, \dots$$

Ejemplo 5.15: Se sabe que en cierto proceso de fabricación uno de cada 100 artículos, en promedio, resulta defectuoso. ¿Cuál es la probabilidad de que el quinto artículo que se inspecciona, en un grupo de 100, sea el primer defectuoso que se encuentra?

Solución: Si utilizamos la distribución geométrica con $x = 5$ y $p = 0.01$, tenemos

$$g(5; 0.01) = (0.01)(0.99)^4 = 0.0096. \quad \text{J}$$

Ejemplo 5.16: En "momentos ajetreados" un conmutador telefónico está muy cerca de su límite de capacidad, por lo que los usuarios tienen dificultad para hacer sus llamadas. Sería interesante saber cuántos intentos serían necesarios para conseguir un enlace telefónico. Suponga que la probabilidad de conseguir un enlace durante un momento ajetreado es $p = 0.05$. Nos interesa conocer la probabilidad de que se necesiten 5 intentos para enlazar con éxito una llamada.

Solución: Si utilizamos la distribución geométrica con $x = 5$ y $p = 0.05$, obtenemos

$$P(X = x) = g(5; 0.05) = (0.05)(0.95)^4 = 0.041. \quad \text{J}$$

Muy a menudo, en aplicaciones que tienen que ver con la distribución geométrica, la media y la varianza son importantes. Se puede ver esto en el ejemplo 5.16, en donde el número *esperado* de llamadas necesario para lograr un enlace es muy importante. A continuación se establecen, sin demostración, la media y la varianza de la distribución geométrica.

Teorema 5.3: La media y la varianza de una variable aleatoria que sigue la distribución geométrica son

$$\mu = \frac{1}{p} \quad \text{y} \quad \sigma^2 = \frac{1-p}{p^2}.$$

Aplicaciones de las distribuciones binomial negativa y geométrica

Las áreas de aplicación de las distribuciones binomial negativa y geométrica serán evidentes cuando nos enfoquemos en los ejemplos de esta sección y en los ejercicios que se dedican a tales distribuciones al final de la sección 5.5. En el caso de la distribución geométrica, el ejemplo 5.16 describe una situación en que los ingenieros o administradores intentan determinar cuán ineficiente es un sistema de conmutación telefónica durante periodos ajetreados. En este caso es evidente que los ensayos que ocurren antes de un éxito representan un costo. Si hay una alta probabilidad de que se requieran varios intentos antes de lograr conectarse, entonces se debería rediseñar el sistema.

Las aplicaciones de la distribución binomial negativa son similares por naturaleza. Supongamos que los intentos son costosos en algún sentido y que *ocurren en secuencia*. La alta probabilidad de que se requiera un número “grande” de intentos para experimentar un número fijo de éxitos no es benéfica ni para el científico ni para el ingeniero. Considere los escenarios de los ejercicios de repaso 5.90 y 5.91. En el ejercicio 5.91 el perforador define cierto nivel de éxitos perforando diferentes sitios en secuencia para encontrar petróleo. Si sólo se han hecho 6 intentos en el momento en que se experimenta el segundo éxito, parecería que las utilidades superan de forma considerable la inversión en que se incurre para la perforación.

5.5 Distribución de Poisson y proceso de Poisson

Los experimentos que producen valores numéricos de una variable aleatoria X , el número de resultados que ocurren durante un intervalo de tiempo determinado o en una región específica, se denominan **experimentos de Poisson**. El intervalo de tiempo puede ser de cualquier duración, como un minuto, un día, una semana, un mes o incluso un año. Por ejemplo, un experimento de Poisson podría generar observaciones para la variable aleatoria X que representa el número de llamadas telefónicas por hora que recibe una oficina, el número de días que una escuela permanece cerrada debido a la nieve durante el invierno o el número de juegos suspendidos debido a la lluvia durante la temporada de béisbol. La región específica podría ser un segmento de recta, una área, un volumen o quizá una pieza de material. En tales casos X podría representar el número de ratas de campo por acre, el número de bacterias en un cultivo dado o el número de errores mecanográficos por página. Un experimento de Poisson se deriva del **proceso de Poisson** y tiene las siguientes propiedades:

Propiedades del proceso de Poisson

1. El número de resultados que ocurren en un intervalo o región específica es independiente del número que ocurre en cualquier otro intervalo de tiempo o región del espacio disjunto. De esta forma vemos que el proceso de Poisson no tiene memoria.
2. La probabilidad de que ocurra un solo resultado durante un intervalo de tiempo muy corto o en una región pequeña es proporcional a la longitud del intervalo o al tamaño de la región, y no depende del número de resultados que ocurren fuera de este intervalo de tiempo o región.
3. La probabilidad de que ocurra más de un resultado en tal intervalo de tiempo corto o que caiga en tal región pequeña es insignificante.

El número X de resultados que ocurren durante un experimento de Poisson se llama **variable aleatoria de Poisson** y su distribución de probabilidad se llama **distribu-**

ción de Poisson. El número medio de resultados se calcula a partir de $\mu = \lambda t$, donde t es el “tiempo”, la “distancia”, el “área” o el “volumen” específicos de interés. Como las probabilidades dependen de λ , denotaremos la tasa de ocurrencia de los resultados con $p(x; \lambda t)$. La derivación de la fórmula para $p(x; \lambda t)$, que se basa en las tres propiedades de un proceso de Poisson que se listaron antes, está fuera del alcance de este texto. La siguiente fórmula se utiliza para calcular probabilidades de Poisson.

Distribución de Poisson La distribución de probabilidad de la variable aleatoria de Poisson X , la cual representa el número de resultados que ocurren en un intervalo de tiempo dado o región específicos y se denota con t , es

$$p(x; \lambda t) = \frac{e^{-\lambda t} (\lambda t)^x}{x!}, \quad x = 0, 1, 2, \dots,$$

donde λ es el número promedio de resultados por unidad de tiempo, distancia, área o volumen y $e = 2.71828\dots$

La tabla A.2 contiene las sumatorias de la probabilidad de Poisson

$$P(r; \lambda t) = \sum_{x=0}^r p(x; \lambda t),$$

para valores selectos de λt que van de 0.1 a 18.0. Ilustramos el uso de esta tabla con los siguientes dos ejemplos.

Ejemplo 5.17: Durante un experimento de laboratorio el número promedio de partículas radiactivas que pasan a través de un contador en un milisegundo es 4. ¿Cuál es la probabilidad de que entren 6 partículas al contador en un milisegundo dado?

Solución: Al usar la distribución de Poisson con $x = 6$ y $\lambda t = 4$, y al remitirnos a la tabla A.2, tenemos que

$$p(6; 4) = \frac{e^{-4} 4^6}{6!} = \sum_{x=0}^6 p(x; 4) - \sum_{x=0}^5 p(x; 4) = 0.8893 - 0.7851 = 0.1042.$$

Ejemplo 5.18: El número promedio de camiones-tanque que llega cada día a cierta ciudad portuaria es 10. Las instalaciones en el puerto pueden alojar a lo sumo 15 camiones-tanque por día. ¿Cuál es la probabilidad de que en un día determinado lleguen más de 15 camiones y se tenga que rechazar algunos?

Solución: Sea X el número de camiones-tanque que llegan cada día. Entonces, usando la tabla A.2, tenemos

$$P(X > 15) = 1 - P(X \leq 15) = 1 - \sum_{x=0}^{15} p(x; 10) = 1 - 0.9513 = 0.0487.$$

Como la distribución binomial, la distribución de Poisson se utiliza para control de calidad, aseguramiento de calidad y muestreo de aceptación. Además, ciertas distribuciones continuas importantes que se usan en la teoría de confiabilidad y en la teoría de colas dependen del proceso de Poisson. Algunas de estas distribuciones se analizan y desarrollan en el capítulo 6. El siguiente teorema acerca de la variable aleatoria de Poisson se presenta en el apéndice A.25.

Teorema 5.4: Tanto la media como la varianza de la distribución de Poisson $p(x; \lambda t)$ son λt .

Naturaleza de la función de probabilidad de Poisson

Al igual que muchas distribuciones discretas y continuas, la forma de la distribución de Poisson se vuelve cada vez más simétrica, incluso con forma de campana, a medida que la media se hace más grande. Una ilustración de esto son las gráficas de la función de probabilidad para $\mu = 0.1$, $\mu = 2$ y finalmente $\mu = 5$ que se muestran en la figura 5.1. Observe cómo se acercan a la simetría cuando μ se vuelve tan grande como 5. Con la distribución binomial ocurre algo parecido, como se ilustrará más adelante en este texto.

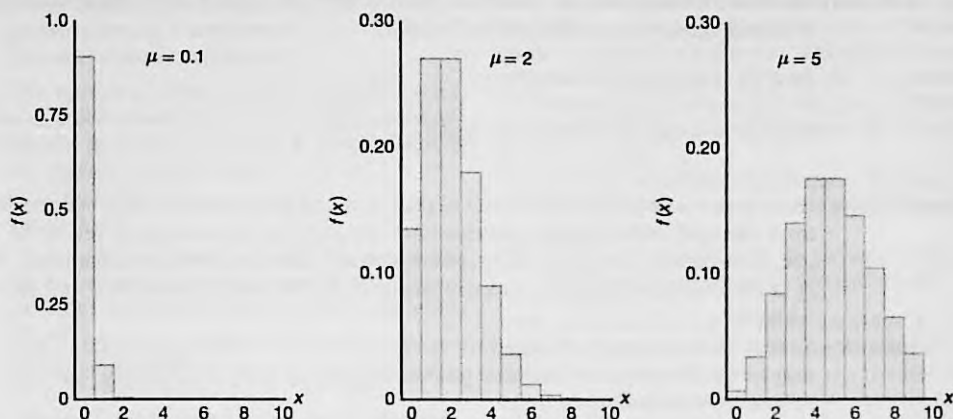


Figura 5.1: Funciones de densidad de Poisson para diferentes medias.

Aproximación de una distribución binomial por medio de una distribución de Poisson

A partir de los tres principios del proceso de Poisson debería ser evidente que la distribución de Poisson se relaciona con la distribución binomial. Aunque la de Poisson por lo general se aplica en problemas de espacio y tiempo, como se ilustra con los ejemplos 5.17 y 5.18, se podría considerar como una forma limitante de la distribución binomial. En el caso de la distribución binomial, si n es bastante grande y p es pequeña, las condiciones comienzan a simular las implicaciones de *espacio* o *tiempo continuos* del proceso de Poisson. La independencia entre las pruebas de Bernoulli en el caso binomial es consistente con la segunda propiedad del proceso de Poisson. Permitir que el parámetro p se acerque a cero se relaciona con la tercera propiedad del proceso de Poisson. De hecho, si n es grande y p es cercana a 0, se puede usar la distribución de Poisson, con $\mu = np$, para aproximar probabilidades binomiales. Si p es cercana a 1, aún podemos utilizar la distribución de Poisson para aproximar probabilidades binomiales intercambiando lo que definimos como éxito y fracaso, por lo tanto, cambiando p a un valor cercano a 0.

Teorema 5.5: Sea X una variable aleatoria binomial con distribución de probabilidad $b(x; n, p)$. Cuando $n \rightarrow \infty$, $p \rightarrow 0$, y $np \xrightarrow{n \rightarrow \infty} \mu$ permanece constante,

$$b(x; n, p) \xrightarrow{n \rightarrow \infty} p(x; \mu).$$

Ejemplo 5.19: En cierta fábrica los accidentes ocurren con muy poca frecuencia. Se sabe que la probabilidad de un accidente en cualquier día dado es de 0.005, y que los accidentes son independientes entre sí.

- a) ¿Cuál es la probabilidad de que en un día de cualquier periodo determinado de 400 días ocurra un accidente?
 b) ¿Cuál es la probabilidad de que ocurra un accidente a lo sumo en tres días de tal periodo?

Solución: Sea X una variable aleatoria binomial con $n = 400$ y $p = 0.005$. Por consiguiente, $np \approx 2$. Si utilizamos la aproximación de Poisson,

$$a) P(X = 1) = e^{-2} 2^1 = 0.271 \text{ y}$$

$$b) P(X \leq 3) = \sum_{x=0}^3 e^{-2} 2^x / x! = 0.857.$$

Ejemplo 5.20: En un proceso de fabricación donde se manufacturan productos de vidrio ocurren defectos o burbujas, lo cual ocasionalmente hace que la pieza ya no se pueda vender. Se sabe que, en promedio, 1 de cada 1000 artículos producidos tiene una o más burbujas. ¿Cuál es la probabilidad de que una muestra aleatoria de 8000 tenga menos de 7 artículos con burbujas?

Solución: Se trata básicamente de un experimento binomial con $n = 8000$ y $p = 0.001$. Como p es muy cercana a cero y n es bastante grande, haremos la aproximación con la distribución de Poisson utilizando

$$\mu = (8000)(0.001) = 8.$$

Por lo tanto, si X representa el número de burbujas, tenemos

$$P(X < 7) = \sum_{x=0}^6 b(x; 8000, 0.001) \approx p(x; 8) = 0.3134.$$

Ejercicios

5.49 La probabilidad de que una persona que vive en cierta ciudad tenga un perro es de 0.3. Calcule la probabilidad de que la décima persona entrevistada al azar en esa ciudad sea la quinta que tiene un perro.

5.50 Calcule la probabilidad de que una persona que lanza una moneda obtenga

- a) la tercera cara en el séptimo lanzamiento;
 b) la primera cara en el cuarto lanzamiento.

5.51 Tres personas lanzan una moneda legal y el disperejo paga los cafés. Si todas las monedas tienen el mismo resultado, se lanzan de nuevo. Calcule la probabilidad de que se necesiten menos de 4 lanzamientos.

5.52 Un científico inyecta a varios ratones, uno a la vez, el virus que produce una enfermedad, hasta que encuentra a 2 que contraen la enfermedad. Si la probabilidad de contraer la enfermedad es de 1/6, ¿cuál es la probabilidad de que tenga que inocular a 8 ratones?

5.53 Un estudio de un inventario determina que, en promedio, el número de veces al día que se solicita un artículo específico en un almacén es 5. ¿Cuál es la probabilidad de que en un día determinado este artículo se pida

- a) más de 5 veces?
 b) ninguna vez?

5.54 De acuerdo con un estudio publicado por un grupo de sociólogos de la Universidad de Massachusetts, Estados Unidos, casi dos terceras partes de los 20 millones de personas que consumen Valium son mujeres. Suponga que esta cifra es una estimación válida y calcule la probabilidad de que en un determinado día la quinta prescripción de Valium que da un médico sea

- a) la primera prescripción de Valium para una mujer;
 b) la tercera prescripción de Valium para una mujer.

5.55 La probabilidad de que una persona que estudia la carrera de piloto privado apruebe el examen escrito para obtener la licencia es de 0.7. Calcule la probabilidad de que cierto estudiante apruebe el examen

- en el tercer intento;
- antes del cuarto intento.

5.56 En cierto crucero ocurren, en promedio, 3 accidentes de tránsito al mes. ¿Cuál es la probabilidad de que en cualquier determinado mes en este crucero

- ocurran exactamente 5 accidentes?
- ocurran menos de 3 accidentes?
- ocurran al menos 2 accidentes?

5.57 Un escritor de libros comete, en promedio, dos errores de procesamiento de texto por página en el primer borrador de su libro. ¿Cuál es la probabilidad de que en la siguiente página cometa

- 4 o más errores?
- ningún error?

5.58 Cierta área del este de Estados Unidos resulta afectada, en promedio, por 6 huracanes al año. Calcule la probabilidad de que para cierto año esta área resulte afectada por

- menos de 4 huracanes;
- cualquier cantidad entre 6 y 8 huracanes.

5.59 Suponga que la probabilidad de que una determinada persona crea un rumor acerca de las transgresiones de cierta actriz famosa es de 0.8. ¿Cuál es la probabilidad de que

- la sexta persona que escuche este rumor sea la cuarta en creerlo?
- la tercera persona que escuche este rumor sea la primera en creerlo?

5.60 Se estima que el número promedio de ratas de campo por acre, en un campo de 5 acres de trigo, es 12. Calcule la probabilidad de que se encuentren menos de 7 ratas de campo

- en un acre dado;
- en 2 de los siguientes 3 acres que se inspeccionen.

5.61 Suponga que, en promedio, una persona en 1000 comete un error numérico al preparar su declaración de impuestos. Si se seleccionan 10,000 formas al azar y se examinan, calcule la probabilidad de que 6, 7 u 8 de las formas contengan un error.

5.62 Se sabe que la probabilidad de que un estudiante de preparatoria no pase la prueba de escoliosis (curvatura de la espina dorsal) es de 0.004. De los siguientes 1875 estudiantes que se revisan en búsqueda de escoliosis, calcule la probabilidad de que

- menos de 5 no pasen la prueba;
- 8, 9 o 10 no pasen la prueba.

5.63 Calcule la media y la varianza de la variable aleatoria X del ejercicio 5.58, que representa el número de huracanes que afectan cada año a cierta área del este de Estados Unidos.

5.64 Calcule la media y la varianza de la variable aleatoria X del ejercicio 5.61, que representa el número de personas, de cada 10,000, que comete un error al preparar su declaración de impuestos.

5.65 Un fabricante de automóviles se preocupa por una falla en el mecanismo de freno de un modelo específico. En raras ocasiones la falla puede causar una catástrofe al manejarlo a alta velocidad. La distribución del número de automóviles por año que experimentará la catástrofe es una variable aleatoria de Poisson con $\lambda = 5$.

- ¿Cuál es la probabilidad de que, a lo sumo, 3 automóviles por año de ese modelo específico sufran una catástrofe?
- ¿Cuál es la probabilidad de que más de un automóvil por año experimente una catástrofe?

5.66 Los cambios en los procedimientos de los aeropuertos requieren una planeación considerable. Los índices de llegadas de los aviones son factores importantes que deben tomarse en cuenta. Suponga que los aviones pequeños llegan a cierto aeropuerto, de acuerdo con un proceso de Poisson, con una frecuencia de 6 por hora. De esta manera, el parámetro de Poisson para las llegadas en un periodo de horas es $\mu = 6t$.

- ¿Cuál es la probabilidad de que lleguen exactamente 4 aviones pequeños durante un periodo de una hora?
- ¿Cuál es la probabilidad de que lleguen al menos 4 durante un periodo de una hora?
- Si definimos un día laboral como de 12 horas, ¿cuál es la probabilidad de que al menos 75 aviones pequeños lleguen durante un día laboral?

5.67 Se supone que el número de clientes que llegan por hora a ciertas instalaciones de servicio automotriz sigue una distribución de Poisson con media $\lambda = 7$.

- Calcule la probabilidad de que lleguen más de 10 clientes en un periodo de dos horas.
- ¿Cuál es el número medio de llegadas durante un periodo de 2 horas?

5.68 Considere el ejercicio 5.62. ¿Cuál es el número promedio de estudiantes que no pasan la prueba?

5.69 La probabilidad de que una persona muera al contraer una infección viral es de 0.001. De los siguientes 4000 infectados con el virus, ¿cuál es el número promedio que morirá?

5.70 Una empresa compra lotes grandes de cierta clase de dispositivo electrónico. Utiliza un método que rechaza el lote completo si en una muestra aleatoria de 100 unidades se encuentran 2 o más unidades defectuosas.

- ¿Cuál es el número promedio de unidades defectuosas que se encuentran en una muestra de 100 unidades si el lote tiene 1% de unidades defectuosas?
- ¿Cuál es la varianza?

5.71 Se sabe que para cierto tipo de alambre de cobre ocurren, en promedio, 1.5 fallas por milímetro. Si se supone que el número de fallas es una variable aleatoria de Poisson, ¿cuál es la probabilidad de que no ocurran fallas en cierta parte de un alambre que tiene 5 milímetros de longitud? ¿Cuál es el número promedio de fallas en alguna parte de un alambre que tiene 5 milímetros de longitud?

5.72 Los baches en ciertas carreteras pueden ser un problema grave y requieren reparación constantemente. Con un tipo específico de terreno y mezcla de concreto la experiencia sugiere que hay, en promedio, 2 baches por milla después de cierta cantidad de uso. Se supone que el proceso de Poisson se aplica a la variable aleatoria "número de baches".

- ¿Cuál es la probabilidad de que no aparezca más de un bache en un tramo de una milla?
- ¿Cuál es la probabilidad de que no aparezcan más de 4 baches en un tramo determinado de 5 millas?

5.73 En ciudades grandes los administradores de los hospitales se preocupan por el flujo de personas en las salas de urgencias. En un hospital específico de una

ciudad grande el personal disponible no puede alojar el flujo de pacientes cuando hay más de 10 casos de emergencia en una hora determinada. Se supone que la llegada de los pacientes sigue un proceso de Poisson y los datos históricos sugieren que, en promedio, llegan 5 emergencias cada hora.

- ¿Cuál es la probabilidad de que en una hora determinada el personal no pueda alojar el flujo de pacientes?
- ¿Cuál es la probabilidad de que, durante un turno de 3 horas, lleguen más de 20 emergencias?

5.74 Se sabe que 3% de las personas a las que se les revisa el equipaje en un aeropuerto lleva objetos cuestionables. ¿Cuál es la probabilidad de que una serie de 15 personas cruce sin problemas antes de que se atrape a una con un objeto cuestionable? ¿Cuál es el número esperado de personas que pasarán antes de que se detenga a una?

5.75 La tecnología cibernética ha generado un ambiente donde los "robots" funcionan con el uso de microprocesadores. La probabilidad de que un robot falle durante cualquier turno de 6 horas es de 0.10. ¿Cuál es la probabilidad de que un robot funcione a lo sumo 5 turnos antes de fallar?

5.76 Se sabe que la tasa de rechazo en las encuestas telefónicas es de aproximadamente 20%. Un reportaje del periódico indica que 50 personas respondieron a una encuesta antes de que una se rehusara a participar.

- Comente acerca de la validez del reportaje. Utilice una probabilidad en su argumento.
- ¿Cuál es el número esperado de personas encuestadas antes de que una se rehúse a responder?

Ejercicios de repaso

5.77 Durante un proceso de producción, cada día se seleccionan al azar 15 unidades de la línea de ensamble para verificar el porcentaje de artículos defectuosos. A partir de información histórica se sabe que la probabilidad de tener una unidad defectuosa es de 0.05. Cada vez que se encuentran dos o más unidades defectuosas en la muestra de 15, el proceso se detiene. Este procedimiento se utiliza para proporcionar una señal en caso de que aumente la probabilidad de unidades defectuosas.

- ¿Cuál es la probabilidad de que en un día determinado se detenga el proceso de producción? (Suponga 5% de unidades defectuosas).
- Suponga que la probabilidad de una unidad defectuosa aumenta a 0.07. ¿Cuál es la probabilidad de que en cualquier día no se detenga el proceso de producción?

5.78 Se considera utilizar una máquina automática de soldadura para un proceso de producción. Antes de comprarla se probará para verificar si tiene éxito en 99% de sus soldaduras. Si no es así, se considerará que no es eficiente. La prueba se llevará a cabo con un prototipo que requiere hacer 100 soldaduras. La máquina se aceptará para la producción sólo si no falla en más de 3 soldaduras.

- ¿Cuál es la probabilidad de que se rechace una buena máquina?
- ¿Cuál es la probabilidad de que se acepte una máquina ineficiente que solde bien el 95% de las veces?

5.79 Una agencia de renta de automóviles en un aeropuerto local tiene 5 Ford, 7 Chevrolet, 4 Dodge, 3 Honda y 4 Toyota disponibles. Si la agencia selecciona al azar 9 de estos automóviles para transportar delega-

dos desde el aeropuerto hasta el centro de convenciones de la ciudad, calcule la probabilidad de que rente 2 Ford, 3 Chevrolet, 1 Dodge, 1 Honda y 2 Toyota.

5.80 En un centro de mantenimiento que recibe llamadas de servicio de acuerdo con un proceso de Poisson entran, en promedio, 2.7 llamadas por minuto. Calcule la probabilidad de que

- no entren más de 4 llamadas en cualquier minuto;
- entren menos de 2 llamadas en cualquier minuto;
- entren más de 10 llamadas en un periodo de 5 minutos.

5.81 Una empresa de electrónica afirma que la proporción de unidades defectuosas de cierto proceso es de 5%. Un comprador sigue el procedimiento estándar de inspeccionar 15 unidades elegidas al azar de un lote grande. En una ocasión específica el comprador encuentra 5 unidades defectuosas.

- ¿Cuál es la probabilidad de que esto ocurra, si es correcta la afirmación de que el 5% de los productos son defectuosos?
- ¿Cómo reaccionaría usted si fuera el comprador?

5.82 Un dispositivo electrónico de conmutación falla ocasionalmente, pero se considera que es satisfactorio si, en promedio, no comete más de 0.20 errores por hora. Se elige un periodo particular de 5 horas para probarlo. Si durante este periodo no ocurre más de un error, se considera que el funcionamiento del dispositivo es satisfactorio.

- ¿Cuál es la probabilidad de que, con base en la prueba, se considere que un dispositivo no funciona satisfactoriamente cuando en realidad sí lo hace? Suponga que se trata de un proceso de Poisson.
- ¿Cuál es la probabilidad de que un dispositivo se considere satisfactorio cuando, de hecho, el número medio de errores que comete es 0.25? De nuevo suponga que se trata de un proceso de Poisson.

5.83 Una empresa por lo general compra lotes grandes de cierta clase de dispositivo electrónico. Utiliza un método que rechaza el lote completo si encuentra 2 o más unidades defectuosas en una muestra aleatoria de 100 unidades.

- ¿Cuál es la probabilidad de que el método rechace un lote que tiene un 1% de unidades defectuosas?
- ¿Cuál es la probabilidad de que acepte un lote que tiene 5% de unidades defectuosas?

5.84 El propietario de una farmacia local sabe que, en promedio, llegan a su farmacia 100 personas por hora.

- Calcule la probabilidad de que en un periodo determinado de 3 minutos nadie entre a la farmacia.
- Calcule la probabilidad de que en un periodo dado de 3 minutos entren más de 5 personas a la farmacia.

5.85 a) Suponga que lanza 4 dados. Calcule la probabilidad de obtener al menos un 1.

b) Suponga que lanza 2 dados 24 veces. Calcule la probabilidad de obtener al menos uno (1, 1), es decir, un "ojos de serpiente".

5.86 Suponga que de 500 billetes de lotería que se venden, 200 le dan a ganar al comprador al menos el costo del billete. Ahora suponga que usted compra 5 billetes. Calcule la probabilidad de ganar al menos el costo de 3 billetes.

5.87 Las imperfecciones en los tableros de circuitos y los microcircuitos de computadora se prestan para un análisis estadístico. Un tipo particular de tablero contiene 200 diodos y la probabilidad de que falle alguno es de 0.03.

- ¿Cuál es el número promedio de fallas en los diodos?
- ¿Cuál es la varianza?
- El tablero funciona si no tiene diodos defectuosos. ¿Cuál es la probabilidad de que un tablero funcione?

5.88 El comprador potencial de un motor particular requiere (entre otras cosas) que éste encienda 10 veces consecutivas. Suponga que la probabilidad de que encienda es de 0.990. Suponga que los resultados de intentos de encendido son independientes.

- ¿Cuál es la probabilidad de que el posible comprador acepte el motor después de sólo 10 encendidos?
- ¿Cuál es la probabilidad de que se tenga que intentar encenderlo 12 veces durante el proceso de aceptación?

5.89 El esquema de aceptación para comprar lotes que contienen un número grande de baterías consiste en probar no más de 75 baterías seleccionadas al azar y rechazar el lote completo si falla una sola batería. Suponga que la probabilidad de encontrar una que falle es de 0.001.

- ¿Cuál es la probabilidad de que se acepte un lote?
- ¿Cuál es la probabilidad de que se rechace un lote en la vigésima prueba?
- ¿Cuál es la probabilidad de que se rechace en 10 o menos pruebas?

5.90 Una empresa que perfora pozos petroleros opera en varios sitios y su éxito o fracaso es independiente de un sitio a otro. Suponga que la probabilidad de éxito en cualquier sitio específico es de 0.25.

- ¿Cuál es la probabilidad de que un perforador burrene 10 sitios y tenga un éxito?
- El perforador se declarará en bancarrota si tiene que perforar 10 veces antes de que ocurra el primer éxito. ¿Cuáles son las perspectivas de bancarrota del perforador?

5.91 Considere la información del ejercicio de repaso 5.90. El perforador cree que "dará en el clavo" si logra el segundo éxito durante o antes del sexto intento. ¿Cuál es la probabilidad de que el perforador "dé en el clavo"?

5.92 Una pareja decide que continuará procreando hijos hasta tener dos hombres. Suponiendo que $P(\text{hombre}) = 0.5$, ¿cuál es la probabilidad de que su segundo niño sea su cuarto hijo?

5.93 Por los investigadores se sabe que una de cada 100 personas es portadora de un gen que lleva a la herencia de cierta enfermedad crónica. En una muestra aleatoria de 1000 individuos, ¿cuál es la probabilidad de que menos de 7 individuos porten el gen? Utilice la aproximación de Poisson. Nuevamente con la aproximación de Poisson, determine cuál es el número promedio aproximado de personas, de cada 1000, que portan el gen.

5.94 Un proceso de fabricación produce piezas para componentes electrónicos. Se supone que la probabilidad de que una pieza salga defectuosa es de 0.01. Durante una prueba de esta suposición se obtiene una muestra al azar de 500 artículos y se encuentran 15 defectuosos.

- ¿Cuál es su respuesta ante la suposición de que 1% de las piezas producidas salen defectuosas? Asegúrese de acompañar su comentario con un cálculo de probabilidad.
- Suponiendo que 1% de las piezas producidas salen con defecto, ¿cuál es la probabilidad de que sólo se encuentren 3 defectuosos?
- Resuelva de nueva cuenta los incisos *a)* y *b)* utilizando la aproximación de Poisson.

5.95 Un proceso de manufactura produce artículos en lotes de 50. Se dispone de planes de muestreo en los cuales los lotes se apartan periódicamente y se someten a cierto tipo de inspección. Por lo general se supone que la proporción de artículos defectuosos que resultan del proceso es muy pequeña. Para la empresa también es importante que los lotes que contengan artículos defectuosos sean un evento raro. El plan actual de inspección consiste en elegir lotes al azar, obtener muestras periódicas de 10 en 50 artículos de un lote y, si ninguno de los muestreados está defectuoso, no se realizan acciones.

- Suponga que se elige un lote al azar y 2 de cada 50 artículos tienen defecto. ¿Cuál es la probabilidad de que al menos uno en la muestra de 10 del lote esté defectuoso?
- A partir de su respuesta en el inciso *a)*, comente sobre la calidad de este plan de muestreo.
- ¿Cuál es el número promedio de artículos defectuosos encontrados por cada 10 artículos de la muestra?

5.96 Considere la situación del ejercicio de repaso 5.95. Se ha determinado que el plan de muestreo debería ser lo suficientemente amplio como para que haya una probabilidad alta, digamos de 0.9, de que si hay tantos como 2 artículos defectuosos en el lote de 50 que se muestrea, al menos uno se encuentre en el muestreo. Con tales restricciones, ¿cuántos de los 50 artículos deberían muestrearse?

5.97 La seguridad nacional requiere que la tecnología de defensa sea capaz de detectar proyectiles o misiles ofensivos. Para que este sistema de defensa sea exitoso, se requieren múltiples pantallas de radar. Suponga que se usarán tres pantallas independientes y que la probabilidad de que cualquiera detecte un misil ofensivo es de 0.8. Es evidente que si ninguna pantalla detecta un misil ofensivo, el sistema no funciona y requiere mejorarse.

- ¿Cuál es la probabilidad de que ninguna de las pantallas detecte un misil ofensivo?
- ¿Cuál es la probabilidad de que sólo una de las pantallas detecte el misil?
- ¿Cuál es la probabilidad de que al menos 2 de las 3 pantallas detecten el misil?

5.98 Suponga que es importante que el sistema general de defensa contra misiles sea lo más perfecto posible.

- Suponga que la calidad de las pantallas es la que se indica en el ejercicio de repaso 5.97. ¿Cuántas se requieren, entonces, para asegurar que la probabilidad de que el misil pase sin ser detectado sea de 0.0001?
- Suponga que se decide utilizar sólo 3 pantallas e intentar mejorar la capacidad de detección de las mismas. ¿Cuál debe ser la eficacia individual de las pantallas (es decir, la probabilidad de detección), para alcanzar la eficacia que se requiere en el inciso *a)*?

5.99 Regrese al ejercicio de repaso 5.95a. Vuelva a calcular la probabilidad usando la distribución binomial. Comente su respuesta.

5.100 En cierto departamento universitario de estadística hay dos vacantes. Cinco personas las solicitan: dos de ellas tienen experiencia con modelos lineales y una tiene experiencia con probabilidad aplicada. Al comité de selección se le indicó elegir a los 2 aspirantes aleatoriamente.

- ¿Cuál es la probabilidad de que los 2 elegidos sean los que tienen experiencia con modelos lineales?
- ¿Cuál es la probabilidad de que, de los 2 elegidos, uno tenga experiencia con modelos lineales y el otro con probabilidad aplicada?

5.101 El fabricante de un triciclo para niños ha recibido quejas porque su producto tiene defecto en los frenos. De acuerdo con el diseño del producto y muchas pruebas preliminares, se determinó que la probabilidad del tipo de defecto reportado era 1 en 10,000 (es decir, de 0.0001). Después de una minuciosa investigación de las quejas se determinó que durante cierto periodo se eligieron aleatoriamente 200 artículos de la producción, de los cuales 5 tuvieron frenos defectuosos.

- a) Comente sobre la afirmación de “uno en 10,000” del fabricante. Utilice un argumento probabilístico. Use la distribución binomial para sus cálculos.
- b) Repita el inciso a utilizando la aproximación de Poisson.

5.102 Proyecto de grupo: Separe la clase en dos grupos aproximadamente del mismo tamaño. Cada uno de los estudiantes del grupo 1 lanzará una moneda 10 veces (n_1) y contará el número de caras resultantes. Cada uno de los estudiantes del grupo 2 lanzará una moneda 40 veces (n_2) y también contará el número de caras obtenidas. Los miembros de cada grupo deben calcular de manera individual la proporción de caras observadas, que es una estimación de p , la probabilidad de obtener una cara. De esta manera, habrá un conjunto de valores de p_1 (del grupo 1) y un conjunto de valores de p_2 (del grupo 2). Todos los valores de p_1 y p_2 son estimaciones de 0.5, que es el valor verdadero de la probabilidad de obtener una cara de una moneda legal.

- a) ¿Cuál conjunto de valores se acerca con mayor consistencia a 0.5, el de p_1 o el de p_2 ? Considere

la demostración del teorema 5.1 de la página 147 con respecto a las estimaciones del parámetro $p = 0.5$. Los valores de p_1 se obtuvieron con $n = n_1 = 10$ y los valores de p_2 se obtuvieron con $n = n_2 = 40$. Si se utiliza la notación de la demostración, las estimaciones están dadas por

$$p_1 = \frac{x_1}{n_1} = \frac{I_1 + \cdots + I_{n_1}}{n_1},$$

donde I_1, \dots, I_{n_1} son ceros y unos y $n_1 = 10$, y

$$p_2 = \frac{x_2}{n_2} = \frac{I_1 + \cdots + I_{n_2}}{n_2},$$

donde I_1, \dots, I_{n_2} son nuevamente ceros y unos y $n_2 = 40$.

- b) Remítase nuevamente al teorema 5.1 y demuestre que

$$E(p_1) = E(p_2) = p = 0.5.$$

- c) Demuestre que $\sigma_{p_1}^2 = \frac{\sigma_{x_1}^2}{n_1}$ es 4 veces el valor de $\sigma_{p_2}^2 = \frac{\sigma_{x_2}^2}{n_2}$. Explique, además, por qué los valores de p_2 del grupo 2 se acercan con mayor consistencia al valor verdadero, $p = 0.5$, que los valores de p_1 del grupo 1.

Aprenderá mucho más sobre la estimación de parámetros a partir del capítulo 9. Ahí pondremos más énfasis en la importancia de la media y la varianza de un estimador de un parámetro.

5.6 Posibles riesgos y errores conceptuales; relación con el material de otros capítulos

Las distribuciones discretas estudiadas en este capítulo ocurren con mucha frecuencia en los escenarios de la ingeniería, así como en los de las ciencias biológicas y físicas. Es evidente que los ejemplos y los ejercicios sugieren esto. Los planes de muestreo industrial y muchas de las decisiones en ingeniería se basan en las distribuciones binomial y de Poisson, así como en la distribución hipergeométrica. Mientras que las distribuciones binomial negativa y geométrica se utilizan en menor grado, también tienen aplicaciones. En específico, una variable aleatoria binomial negativa se puede ver como una mezcla de variables aleatorias gamma y de Poisson (la distribución gamma se estudiará en el capítulo 6).

A pesar de las múltiples aplicaciones que estas distribuciones tienen en la vida real, podrían utilizarse de manera incorrecta, a menos que el científico sea prudente y cuidadoso. Desde luego, cualquier cálculo de probabilidad para las distribuciones que se estudiaron en este capítulo se realiza bajo el supuesto de que se conoce el valor del parámetro. Las aplicaciones en el mundo real a menudo resultan en un valor del parámetro que se puede “desplazar” debido a factores que son difíciles de controlar en el proceso,

o debido a intervenciones en el proceso que no se han tomado en cuenta. Por ejemplo, en el ejercicio de repaso 5.77 se utilizó "información histórica"; sin embargo, ¿el proceso actual es el mismo que aquel en que se recabaron los datos históricos? El uso de la distribución de Poisson tiene incluso más posibilidades de enfrentar esta dificultad. Por ejemplo, en el ejercicio de repaso 5.80 las preguntas de los incisos *a*, *b* y *c* se basan en el uso de $\mu = 2.7$ llamadas por minuto. Con base en los registros históricos éste es el número de llamadas que se reciben "en promedio". Pero en ésta y muchas otras aplicaciones de la distribución de Poisson hay momentos desocupados y momentos ajetreados, de manera que se espera que haya momentos en que las condiciones para el proceso de Poisson parezcan cumplirse, cuando en realidad no lo hacen. Por consiguiente, los cálculos de probabilidad podrían ser incorrectos. En el caso de la distribución binomial, la condición que podría fallar en ciertas aplicaciones (además de la falta de constancia de p) es la suposición de independencia, estipulando que los experimentos de Bernoulli son independientes.

Una de las aplicaciones incorrectas más célebres de la distribución binomial ocurrió en la temporada de béisbol de 1961, cuando Mickey Mantle y Roger Maris se enfrascaron en una batalla amistosa por romper el récord de todos los tiempos de 60 jonrones establecido por Babe Ruth. Un famoso artículo de una revista predijo, con base en la teoría de la probabilidad, que Mantle rompería el récord. La predicción estaba fundamentada en un cálculo de probabilidad en el que se utilizó la distribución binomial. El error clásico cometido fue la estimación del parámetro p (uno para cada jugador) con base en la frecuencia histórica relativa de jonrones a lo largo de la carrera de los 2 jugadores. Maris, a diferencia de Mantle, no había sido un jonronero prodigioso antes de 1961, de manera que su estimado de p fue bastante bajo. Como resultado de esto se determinó que Mantle tenía más probabilidades que Maris de romper el récord, pero quien logró romperlo al final fue este último.

Capítulo 6

Algunas distribuciones continuas de probabilidad

6.1 Distribución uniforme continua

Una de las distribuciones continuas más simples de la estadística es la **distribución uniforme continua**. Esta distribución se caracteriza por una función de densidad que es “plana”, por lo cual la probabilidad es uniforme en un intervalo cerrado, digamos $[A, B]$. Aunque las aplicaciones de la distribución uniforme continua no son tan abundantes como las de otras distribuciones que se presentan en este capítulo, es apropiado para el principiante que comience esta introducción a las distribuciones continuas con la distribución uniforme.

Distribución uniforme La función de densidad de la variable aleatoria uniforme continua X en el intervalo $[A, B]$ es

$$f(x; A, B) = \begin{cases} \frac{1}{B-A}, & A \leq x \leq B, \\ 0, & \text{en otro caso.} \end{cases}$$

La función de densidad forma un rectángulo con base $B - A$ y **altura constante** $\frac{1}{B-A}$. Como resultado, la distribución uniforme a menudo se conoce como **distribución rectangular**. Sin embargo, observe que el intervalo no siempre es cerrado: $[A, B]$; también puede ser (A, B) . En la figura 6.1 se muestra la función de densidad para una variable aleatoria uniforme en el intervalo $[1, 3]$.

Resulta sencillo calcular las probabilidades para la distribución uniforme debido a la naturaleza simple de la función de densidad. Sin embargo, observe que la aplicación de esta distribución se basa en el supuesto de que la probabilidad de caer en un intervalo de longitud fija dentro de $[A, B]$ es constante.

Ejemplo 6.1: Suponga que el tiempo máximo que se puede reservar una sala de conferencias grande de cierta empresa son cuatro horas. Con mucha frecuencia tienen conferencias extensas y breves. De hecho, se puede suponer que la duración X de una conferencia tiene una distribución uniforme en el intervalo $[0, 4]$.

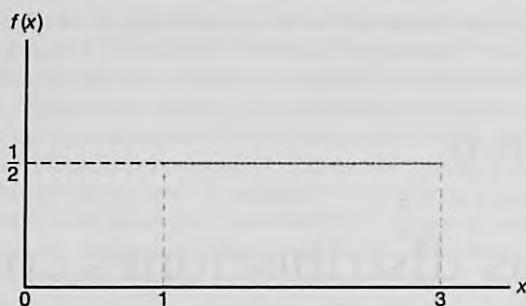


Figura 6.1: Función de densidad para una variable aleatoria en el intervalo $[1, 3]$.

- a) ¿Cuál es la función de densidad de probabilidad?
 b) ¿Cuál es la probabilidad de que cualquier conferencia determinada dure al menos 3 horas?

Solución: a) La función de densidad apropiada para la variable aleatoria X distribuida uniformemente en esta situación es

$$f(x) = \begin{cases} \frac{1}{4}, & 0 \leq x \leq 4, \\ 0, & \text{en otro caso.} \end{cases}$$

b) $P[X \geq 3] = \int_3^4 \frac{1}{4} dx = \frac{1}{4}$.

Teorema 6.1: La media y la varianza de la distribución uniforme son

$$\mu = \frac{A+B}{2} \text{ y } \sigma^2 = \frac{(B-A)^2}{12}.$$

Las demostraciones de los teoremas se dejan al lector. Véase el ejercicio 6.1 de la página 185.

6.2 Distribución normal

La distribución de probabilidad continua más importante en todo el campo de la estadística es la **distribución normal**. Su gráfica, denominada **curva normal**, es la curva con forma de campana de la figura 6.2, la cual describe de manera aproximada muchos fenómenos que ocurren en la naturaleza, la industria y la investigación. Por ejemplo, las mediciones físicas en áreas como los experimentos meteorológicos, estudios de la precipitación pluvial y mediciones de partes fabricadas a menudo se explican más que adecuadamente con una distribución normal. Además, los errores en las mediciones científicas se aproximan muy bien mediante una distribución normal. En 1733, Abraham DeMoivre desarrolló la ecuación matemática de la curva normal, la cual sentó las bases sobre las que descansa gran parte de la teoría de la estadística inductiva. La distribución normal a menudo se denomina **distribución gaussiana** en honor de Karl Friedrich Gauss (1777-1855), quien también derivó su ecuación a partir de un estudio de errores en mediciones repetidas de la misma cantidad.

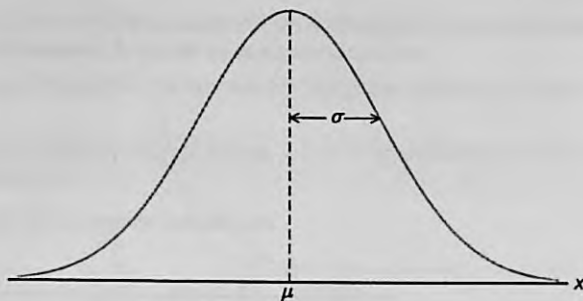


Figura 6.2: La curva normal.

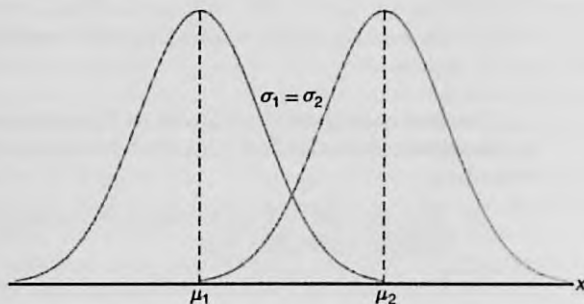
Una variable aleatoria continua X que tiene la distribución en forma de campana de la figura 6.2 se denomina **variable aleatoria normal**. La ecuación matemática para la distribución de probabilidad de la variable normal depende de los dos parámetros μ y σ , su media y su desviación estándar, respectivamente. Por ello, denotamos los valores de la densidad de X por $n(x; \mu, \sigma)$.

Distribución normal La densidad de la variable aleatoria normal X , con media μ y varianza σ^2 , es

$$n(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}, \quad -\infty < x < \infty,$$

donde $\pi = 3.14159\dots$ y $e = 2.71828\dots$

Una vez que se especifican μ y σ , la curva normal queda determinada por completo. Por ejemplo, si $\mu = 50$ y $\sigma = 5$, entonces se pueden calcular las ordenadas $n(x; 50, 5)$ para diferentes valores de x y dibujar la curva. En la figura 6.3 aparecen dos curvas normales que tienen la misma desviación estándar pero diferentes medias. Las dos curvas son idénticas en forma, pero están centradas en diferentes posiciones a lo largo del eje horizontal.

Figura 6.3: Curvas normales con $\mu_1 < \mu_2$ y $\sigma_1 = \sigma_2$.

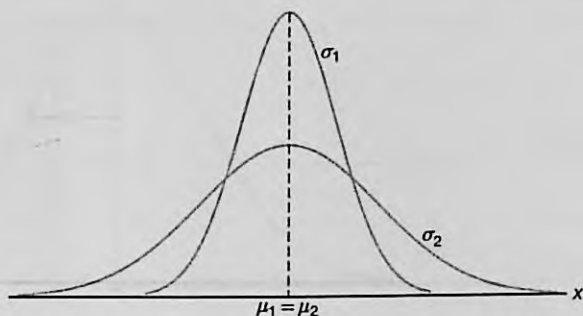


Figura 6.4: Curvas normales con $\mu_1 = \mu_2$ y $\sigma_1 < \sigma_2$.

En la figura 6.4 se muestran dos curvas normales con la misma media pero con desviaciones estándar diferentes. Aquí se observa que las dos curvas están centradas exactamente en la misma posición sobre el eje horizontal; sin embargo, la curva con la mayor desviación estándar es más baja y más extendida. Recuerde que el área bajo una curva de probabilidad debe ser igual a 1 y, por lo tanto, cuanto más variable sea el conjunto de observaciones, más baja y más ancha será la curva correspondiente.

La figura 6.5 muestra dos curvas normales que tienen diferentes medias y diferentes desviaciones estándar. Evidentemente, están centradas en posiciones diferentes sobre el eje horizontal y sus formas reflejan los dos valores diferentes de σ .

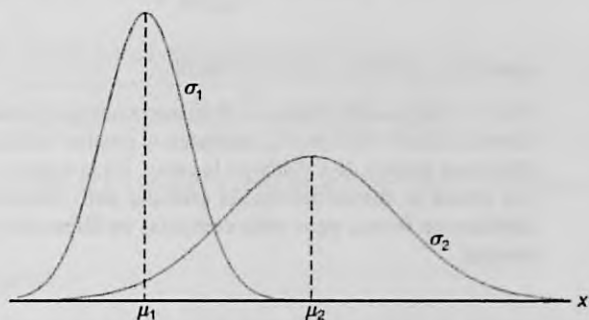


Figura 6.5: Curvas normales con $\mu_1 < \mu_2$ y $\sigma_1 < \sigma_2$.

Con base en lo que observamos en las figuras 6.2 a 6.5, y en el examen de la primera y la segunda derivadas de $n(x; \mu, \sigma)$, listamos las siguientes propiedades de la curva normal:

1. La moda, que es el punto sobre el eje horizontal donde la curva tiene su punto máximo, ocurre en $x = \mu$.
2. La curva es simétrica alrededor de un eje vertical a través de la media μ .
3. La curva tiene sus puntos de inflexión en $x = \mu \pm \sigma$, es cóncava hacia abajo si $\mu - \sigma < X < \mu + \sigma$, y es cóncava hacia arriba en otro caso.

4. La curva normal se aproxima al eje horizontal de manera asintótica, conforme nos alejamos de la media en cualquier dirección.
5. El área total bajo la curva y sobre el eje horizontal es igual a uno.

Teorema 6.2: La media y la varianza de $n(x; \mu, \sigma)$ son μ y σ^2 , respectivamente. Por lo tanto, la desviación estándar es σ .

Prueba: Para evaluar la media primero calculamos

$$E(X - \mu) = \int_{-\infty}^{\infty} \frac{x - \mu}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx.$$

Al establecer que $z = (x - \mu)/\sigma$ y $dx = \sigma dz$, obtenemos

$$E(X - \mu) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z e^{-\frac{1}{2}z^2} dz = 0,$$

dado que la integral anterior es una función impar de z . Al aplicar el teorema 4.5 de la página 128 concluimos que

$$E(X) = \mu$$

La varianza de la distribución normal es dada por

$$E[(X - \mu)^2] = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (x - \mu)^2 e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx.$$

De nuevo, al establecer que $z = (x - \mu)/\sigma$ y $dx = \sigma dz$, obtenemos

$$E[(X - \mu)^2] = \frac{\sigma^2}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z^2 e^{-\frac{1}{2}z^2} dz.$$

Al integrar por partes con $u = z$ y $dv = z e^{-z^2/2} dz$ de modo que $du = dz$ y $v = -e^{-z^2/2}$, encontramos que

$$E[(X - \mu)^2] = \frac{\sigma^2}{\sqrt{2\pi}} \left(-ze^{-z^2/2} \Big|_{-\infty}^{\infty} + \int_{-\infty}^{\infty} e^{-z^2/2} dz \right) = \sigma^2(0 + 1) = \sigma^2. \quad \blacksquare$$

Muchas variables aleatorias tienen distribuciones de probabilidad que se pueden describir de forma adecuada mediante la curva normal, una vez que se especifiquen μ y σ^2 . En este capítulo supondremos que se conocen estos dos parámetros, quizás a partir de investigaciones anteriores. Más adelante haremos inferencias estadísticas cuando se desconozcan μ y σ^2 y se estimen a partir de los datos experimentales disponibles.

Anteriormente señalamos el papel que desempeña la distribución normal como una aproximación razonable de variables científicas en experimentos de la vida real. Hay otras aplicaciones de la distribución normal que el lector apreciará a medida que avance en el estudio de este libro. La distribución normal tiene muchas aplicaciones como *distribución limitante*. En ciertas condiciones, la distribución normal ofrece una buena aproximación continua a las distribuciones binomial e hipergeométrica. El caso de la aproximación a la distribución binomial se examina en la sección 6.5. En el capítulo 8 el lector aprenderá acerca de las **distribuciones muestrales**. Resulta que la distribución limitante de promedios muestrales es normal, lo que brinda una base amplia para la

inferencia estadística, que es muy valiosa para el analista de datos interesado en la estimación y prueba de hipótesis. Las teorías de áreas importantes como el análisis de varianzas (capítulos 13, 14 y 15) y el control de calidad (capítulo 17) se basan en suposiciones que utilizan la distribución normal.

En la sección 6.3 se ofrecen ejemplos para demostrar cómo se utilizan las tablas de la distribución normal. En la sección 6.4 continúan los ejemplos de aplicaciones de la distribución normal.

6.3 Áreas bajo la curva normal

La curva de cualquier distribución continua de probabilidad o función de densidad se construye de manera que el área bajo la curva limitada por las dos ordenadas $x = x_1$ y $x = x_2$ sea igual a la probabilidad de que la variable aleatoria X tome un valor entre $x = x_1$ y $x = x_2$. Por consiguiente, para la curva normal de la figura 6.6,

$$P(x_1 < X < x_2) = \int_{x_1}^{x_2} n(x; \mu, \sigma) dx = \frac{1}{\sqrt{2\pi}\sigma} \int_{x_1}^{x_2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx,$$

es representada por el área de la región sombreada.

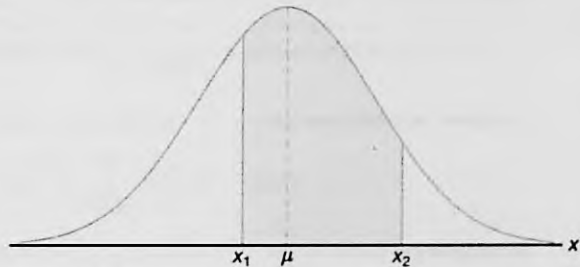


Figura 6.6: $P(x_1 < X < x_2) = \text{área de la región sombreada}$.

En las figuras 6.3, 6.4 y 6.5 vimos cómo la curva normal depende de la media y de la desviación estándar de la distribución que se está estudiando. El área bajo la curva entre cualesquiera dos ordenadas también debe depender de los valores μ y σ . Esto es evidente en la figura 6.7, donde sombreamos las regiones que corresponden a $P(x_1 < X < x_2)$ para dos curvas con medias y varianzas diferentes. $P(x_1 < X < x_2)$, donde X es la variable aleatoria que describe la distribución A , se indica por el área sombreada más oscura debajo de la curva de A . Si X es la variable aleatoria que describe la distribución B , entonces $P(x_1 < X < x_2)$ es dada por toda la región sombreada. Evidentemente, las dos regiones sombreadas tienen tamaños diferentes; por lo tanto, la probabilidad asociada con cada distribución será diferente para los dos valores dados de X .

Existen muchos tipos de programas estadísticos que sirven para calcular el área bajo la curva normal. La dificultad que se enfrenta al resolver las integrales de funciones de densidad normal exige tabular las áreas de la curva normal para una referencia rápida. Sin embargo, sería inútil tratar de establecer tablas separadas para cada posible valor de μ y σ . Por fortuna, podemos transformar todas las observaciones de cualquier variable

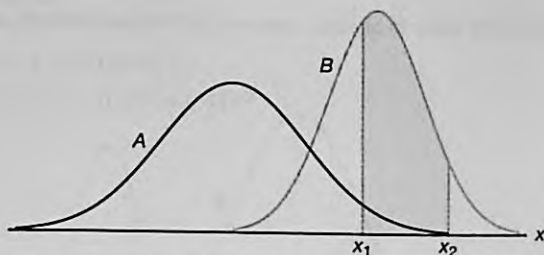


Figura 6.7: $P(x_1 < X < x_2)$ para diferentes curvas normales.

aleatoria normal X en un nuevo conjunto de observaciones de una variable aleatoria normal Z con media 0 y varianza 1. Esto se puede realizar mediante la transformación

$$Z = \frac{X - \mu}{\sigma}.$$

Siempre que X tome un valor x , el valor correspondiente de Z es dado por $z = (x - \mu)/\sigma$. Por lo tanto, si X cae entre los valores $x = x_1$ y $x = x_2$, la variable aleatoria Z caerá entre los valores correspondientes $z_1 = (x_1 - \mu)/\sigma$ y $z_2 = (x_2 - \mu)/\sigma$. En consecuencia, podemos escribir

$$\begin{aligned} P(x_1 < X < x_2) &= \frac{1}{\sqrt{2\pi}\sigma} \int_{x_1}^{x_2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx = \frac{1}{\sqrt{2\pi}} \int_{z_1}^{z_2} e^{-\frac{1}{2}z^2} dz \\ &= \int_{z_1}^{z_2} n(z; 0, 1) dz = P(z_1 < Z < z_2), \end{aligned}$$

donde Z se considera una variable aleatoria normal con media 0 y varianza 1.

Definición 6.1: La distribución de una variable aleatoria normal con media 0 y varianza 1 se llama **distribución normal estándar**.

Las distribuciones original y transformada se ilustran en la figura 6.8. Como todos los valores de X que caen entre x_1 y x_2 tienen valores z correspondientes entre z_1 y z_2 , el área bajo la curva X entre las ordenadas $x = x_1$ y $x = x_2$ de la figura 6.8 es igual al área bajo la curva Z entre las ordenadas transformadas $z = z_1$ y $z = z_2$.

Ahora hemos reducido el número requerido de tablas de áreas de curva normal a una, la de la distribución normal estándar. La tabla A.3 indica el área bajo la curva normal estándar que corresponde a $P(Z < z)$ para valores de z que van de -3.49 a 3.49 . Para ilustrar el uso de esta tabla calculemos la probabilidad de que Z sea menor que 1.74 . Primero, localizamos un valor de z igual a 1.7 en la columna izquierda, después nos movemos a lo largo del renglón hasta la columna bajo 0.04 , donde leemos 0.9591 . Por lo tanto, $P(Z < 1.74) = 0.9591$. Para calcular un valor z que corresponda a una probabilidad dada se invierte el proceso. Por ejemplo, se observa que el valor z que deja un área de 0.2148 bajo la curva a la izquierda de z es -0.79 .

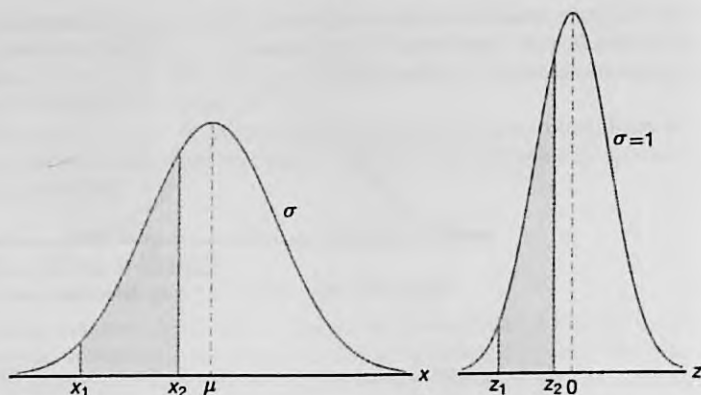


Figura 6.8: Distribuciones normales original y transformada.

Ejemplo 6.2: Dada una distribución normal estándar, calcule el área bajo la curva que se localiza

- a la derecha de $z = 1.84$, y
- entre $z = -1.97$ y $z = 0.86$.

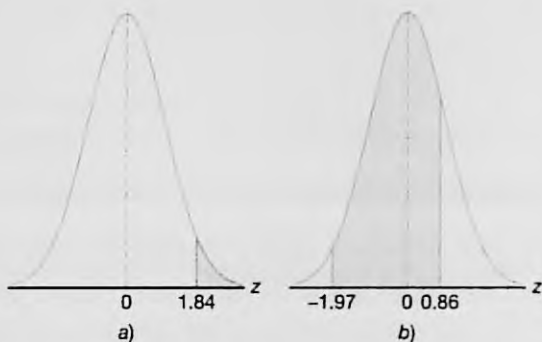


Figura 6.9: Áreas para el ejemplo 6.2.

Solución: Véase la figura 6.9 para las áreas específicas.

- El área en la figura 6.9a a la derecha de $z = 1.84$ es igual a 1 menos el área en la tabla A.3 a la izquierda de $z = 1.84$, a saber, $1 - 0.9671 = 0.0329$.
- El área en la figura 6.9b entre $z = -1.97$ y $z = 0.86$ es igual al área a la izquierda de $z = 0.86$ menos el área a la izquierda de $z = -1.97$. A partir de la tabla A.3 encontramos que el área que se desea es $0.8051 - 0.0244 = 0.7807$. J

Ejemplo 6.3: Dada una distribución normal estándar, calcule el valor de k tal que

a) $P(Z > k) = 0.3015$, y

b) $P(k < Z < -0.18) = 0.4197$.

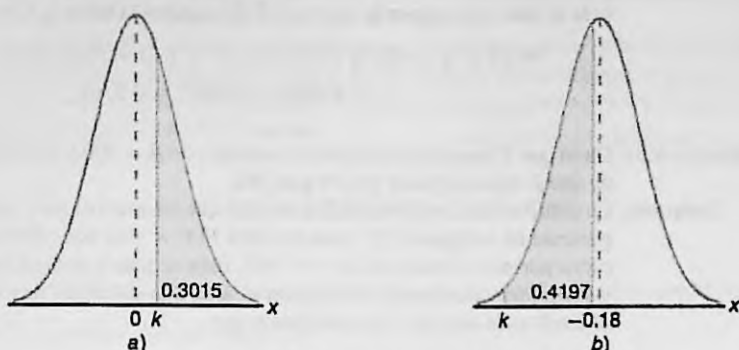


Figura 6.10: Áreas para el ejemplo 6.3.

Solución: La distribución y las áreas deseadas se muestran en la figura 6.10.

- a) En la figura 6.10a vemos que el valor k que deja un área de 0.3015 a la derecha debe dejar entonces un área de 0.6985 a la izquierda. De la tabla A.3 se sigue que $k = 0.52$.
- b) En la tabla A.3 observamos el área total a la izquierda de -0.18 es igual a 0.4286. En la figura 6.10b vemos que el área entre k y -0.18 es 0.4197, de manera que el área a la izquierda de k debe ser $0.4286 - 0.4197 = 0.0089$. Por lo tanto, a partir de la tabla A.3 tenemos $k = -2.37$. ■

Ejemplo 6.4: Dada una variable aleatoria X que tiene una distribución normal con $\mu = 50$ y $\sigma = 10$, calcule la probabilidad de que X tome un valor entre 45 y 62.

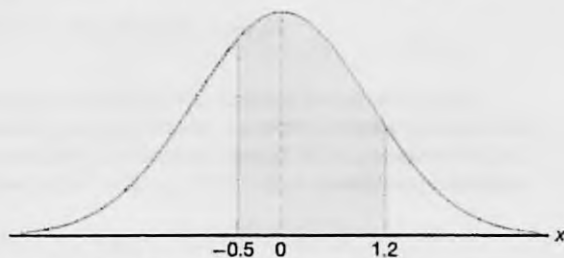


Figura 6.11: Área para el ejemplo 6.4.

Solución: Los valores z que corresponden a $x_1 = 45$ y $x_2 = 62$ son

$$z_1 = \frac{45 - 50}{10} = -0.5 \text{ y } z_2 = \frac{62 - 50}{10} = 1.2.$$

Por lo tanto,

$$P(45 < X < 62) = P(-0.5 < Z < 1.2).$$

$P(-0.5 < Z < 1.2)$ se muestra mediante el área de la región sombreada de la figura 6.11. Esta área se puede calcular restando el área a la izquierda de la ordenada $z = -0.5$ de toda el área a la izquierda de $z = 1.2$. Si usamos la tabla A.3, tenemos

$$\begin{aligned} P(45 < X < 62) &= P(-0.5 < Z < 1.2) = P(Z < 1.2) - P(Z < -0.5) \\ &= 0.8849 - 0.3085 = 0.5764. \end{aligned}$$

Ejemplo 6.5: Dado que X tiene una distribución normal con $\mu = 300$ y $\sigma = 50$, calcule la probabilidad de que X tome un valor mayor que 362.

Solución: La distribución de probabilidad normal que muestra el área sombreada que se desea se presenta en la figura 6.12. Para calcular $P(X > 362)$ necesitamos evaluar el área bajo la curva normal a la derecha de $x = 362$. Esto se puede realizar transformando $x = 362$ al valor z correspondiente, obteniendo el área a la izquierda de z de la tabla A.3 y después restando esta área de 1. Encontramos que

$$z = \frac{362 - 300}{50} = 1.24.$$

De ahí,

$$P(X > 362) = P(Z > 1.24) = 1 - P(Z < 1.24) = 1 - 0.8925 = 0.1075. \quad \text{J}$$

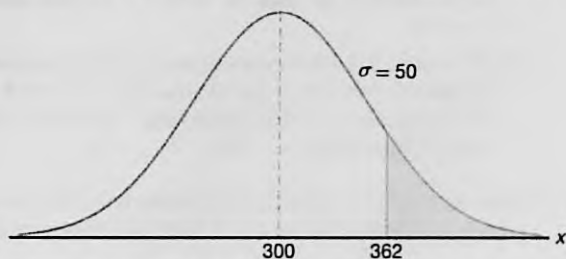


Figura 6.12: Área para el ejemplo 6.5.

De acuerdo con el teorema de Chebyshev en la página 137, la probabilidad de que una variable aleatoria tome un valor dentro de 2 desviaciones estándar de la media es de por lo menos $3/4$. Si la variable aleatoria tiene una distribución normal, los valores z que corresponden a $x_1 = \mu - 2\sigma$ y $x_2 = \mu + 2\sigma$ se calculan fácilmente y son

$$z_1 = \frac{(\mu - 2\sigma) - \mu}{\sigma} = -2 \text{ y } z_2 = \frac{(\mu + 2\sigma) - \mu}{\sigma} = 2.$$

De ahí,

$$\begin{aligned} P(\mu - 2\sigma < X < \mu + 2\sigma) &= P(-2 < Z < 2) = P(Z < 2) - P(Z < -2) \\ &= 0.9772 - 0.0228 = 0.9544, \end{aligned}$$

que es una afirmación mucho más firme que la que se establece mediante el teorema de Chebyshev.

Uso de la curva normal a la inversa

En ocasiones se nos pide calcular el valor de z que corresponde a una probabilidad específica que cae entre los valores que se listan en la tabla A.3 (véase el ejemplo 6.6). Por conveniencia, siempre elegiremos el valor z que corresponde a la probabilidad tabular que está más cerca de la probabilidad que se especifica.

Los dos ejemplos anteriores se resolvieron al ir primero de un valor de x a un valor z y después calcular el área que se desea. En el ejemplo 6.6 invertimos el proceso y comenzamos con un área o probabilidad conocida, calculamos el valor z y después determinamos x reacomodando la fórmula

$$z = \frac{x - \mu}{\sigma} \text{ para obtener } x = \sigma z + \mu.$$

Ejemplo 6.6: Dada una distribución normal con $\mu = 40$ y $\sigma = 6$, calcule el valor de x que tiene

- 45% del área a la izquierda, y
- 14% del área a la derecha.

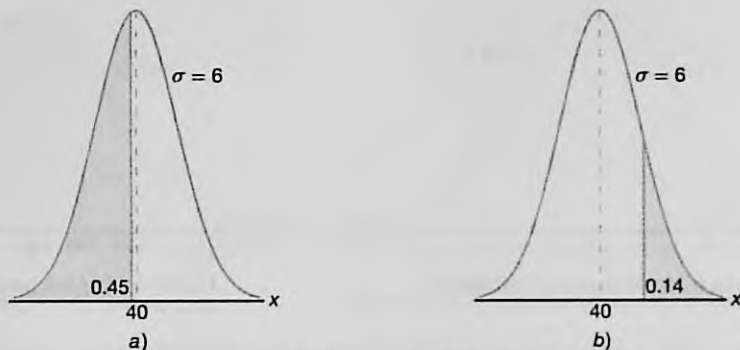


Figura 6.13: Áreas para el ejemplo 6.6.

Solución: a) En la figura 6.13a se sombrea un área de 0.45 a la izquierda del valor x deseado. Necesitamos un valor z que deje un área de 0.45 a la izquierda. En la tabla A.3 encontramos $P(Z < -0.13) = 0.45$, es decir, que el valor z que se desea es -0.13 . Por lo tanto,

$$x = (6)(-0.13) + 40 = 39.22.$$

- b) En la figura 6.13b sombreamos un área igual a 0.14 a la derecha del valor x deseado. Esta vez necesitamos un valor z que deje 0.14 del área a la derecha y, por lo tanto, un área de 0.86 a la izquierda. De nuevo, a partir de la tabla A.3 encontramos $P(Z < 1.08) = 0.86$, así que el valor z deseado es 1.08 y

$$x = (6)(1.08) + 40 = 46.48.$$

6.4 Aplicaciones de la distribución normal

En los siguientes ejemplos se abordan algunos de los muchos problemas en los que se puede aplicar la distribución normal. El uso de la curva normal para aproximar probabilidades binomiales se estudia en la sección 6.5.

Ejemplo 6.7: Cierta tipo de batería de almacenamiento dura, en promedio, 3.0 años, con una desviación estándar de 0.5 años. Suponga que la duración de la batería se distribuye normalmente y calcule la probabilidad de que una batería determinada dure menos de 2.3 años.

Solución: Empiece construyendo un diagrama como el de la figura 6.14, que muestra la distribución dada de la duración de las baterías y el área deseada. Para calcular la $P(X < 2.3)$ necesitamos evaluar el área bajo la curva normal a la izquierda de 2.3. Esto se logra calculando el área a la izquierda del valor z correspondiente. De donde encontramos que

$$z = \frac{2.3 - 3}{0.5} = -1.4,$$

y entonces, usando la tabla A.3, tenemos

$$P(X < 2.3) = P(Z < -1.4) = 0.0808. \quad \lrcorner$$

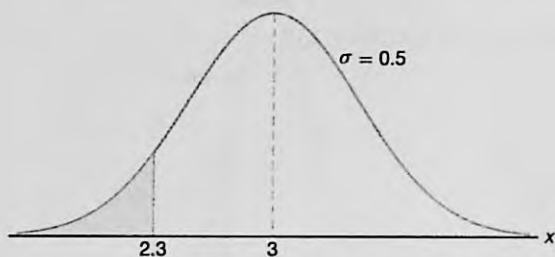


Figura 6.14: Área para el ejemplo 6.7.

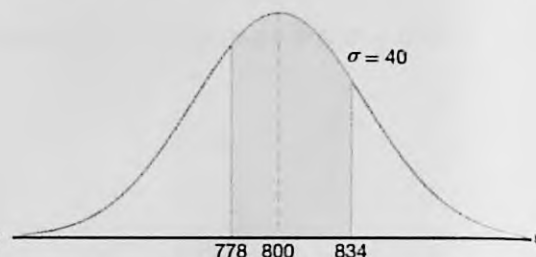


Figura 6.15: Área para el ejemplo 6.8.

Ejemplo 6.8: Una empresa de material eléctrico fabrica bombillas de luz cuya duración, antes de quemarse, se distribuye normalmente con una media igual a 800 horas y una desviación estándar de 40 horas. Calcule la probabilidad de que una bombilla se quemara entre 778 y 834 horas.

Solución: La distribución de vida de las bombillas se ilustra en la figura 6.15. Los valores z que corresponden a $x_1 = 778$ y $x_2 = 834$ son

$$z_1 = \frac{778 - 800}{40} = -0.55 \text{ y } z_2 = \frac{834 - 800}{40} = 0.85.$$

Por lo tanto,

$$\begin{aligned} P(778 < X < 834) &= P(-0.55 < Z < 0.85) = P(Z < 0.85) - P(Z < -0.55) \\ &= 0.8023 - 0.2912 = 0.5111. \quad \lrcorner \end{aligned}$$

Ejemplo 6.9: En un proceso industrial el diámetro de un cojinete de bolas es una medida importante. El comprador establece que las especificaciones en el diámetro sean 3.0 ± 0.01 cm. Esto

implica que no se aceptará ninguna parte que no cumpla estas especificaciones. Se sabe que en el proceso el diámetro de un cojinete tiene una distribución normal con media $\mu = 3.0$ y una desviación estándar $\sigma = 0.005$. En promedio, ¿cuántos de los cojinetes fabricados se descartarán?

Solución: La distribución de los diámetros se ilustra en la figura 6.16. Los valores que corresponden a los límites especificados son $x_1 = 2.99$ y $x_2 = 3.01$. Los valores z correspondientes son

$$z_1 = \frac{2.99 - 3.0}{0.005} = -2.0 \text{ y } z_2 = \frac{3.01 - 3.0}{0.005} = +2.0.$$

Por lo tanto,

$$P(2.99 < X < 3.01) = P(-2.0 < Z < 2.0).$$

A partir de la tabla A.3, $P(Z < -2.0) = 0.0228$. Debido a la simetría de la distribución normal, encontramos que

$$P(Z < -2.0) + P(Z > 2.0) = 2(0.0228) = 0.0456.$$

Como resultado se anticipa que, en promedio, se descartarán 4.56% de los cojinetes fabricados. ▀

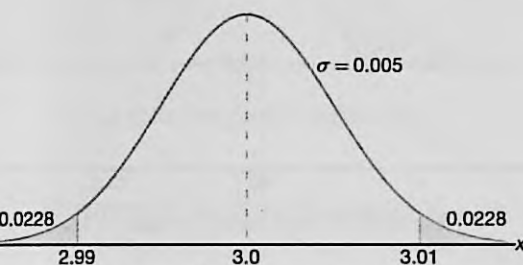


Figura 6.16: Área para el ejemplo 6.9.

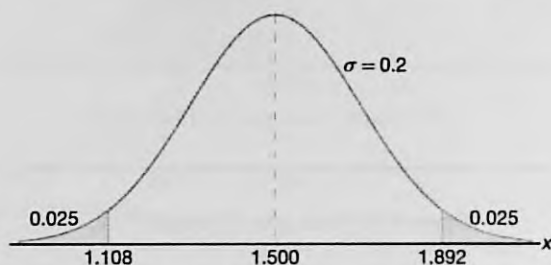


Figura 6.17: Especificaciones para el ejemplo 6.10.

Ejemplo 6.10: Se utilizan medidores para rechazar todos los componentes en los que cierta dimensión no esté dentro de la especificación $1.50 \pm d$. Se sabe que esta medida se distribuye normalmente con una media de 1.50 y una desviación estándar de 0.2. Determine el valor d tal que las especificaciones “cubran” 95% de las mediciones.

Solución: A partir de la tabla A.3 sabemos que

$$P(-1.96 < Z < 1.96) = 0.95.$$

Por lo tanto,

$$1.96 = \frac{(1.50 + d) - 1.50}{0.2},$$

de la que obtenemos

$$d = (0.2)(1.96) = 0.392.$$

En la figura 6.17 se muestra una ilustración de las especificaciones. ▀

Ejemplo 6.11: Cierta máquina fabrica resistencias eléctricas que tienen una resistencia media de 40 ohms y una desviación estándar de 2 ohms. Si se supone que la resistencia sigue una distribución normal y que se puede medir con cualquier grado de precisión, ¿qué porcentaje de resistencias tendrán una resistencia que exceda 43 ohms?

Solución: Se obtiene un porcentaje multiplicando la frecuencia relativa por 100%. Como la frecuencia relativa para un intervalo es igual a la probabilidad de caer en el intervalo, debemos calcular el área a la derecha de $x = 43$ en la figura 6.18. Esto se puede hacer transformando $x = 43$ al valor z correspondiente, con lo cual se obtiene el área a la izquierda de z de la tabla A.3, y después se resta esta área de 1. Encontramos que

$$z = \frac{43 - 40}{2} = 1.5.$$

Por lo tanto,

$$P(X > 43) = P(Z > 1.5) = 1 - P(Z < 1.5) = 1 - 0.9332 = 0.0668.$$

Así, 6.68% de las resistencias tendrán una resistencia que exceda 43 ohms. ┘

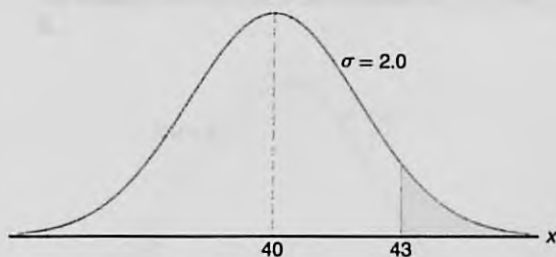


Figura 6.18: Área para el ejemplo 6.11.

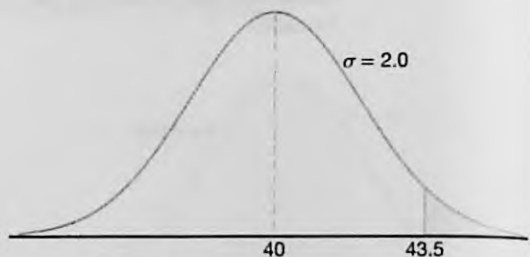


Figura 6.19: Área para el ejemplo 6.12.

Ejemplo 6.12: Calcule el porcentaje de resistencias que excedan 43 ohms para el ejemplo 6.11 si la resistencia se mide al ohm más cercano.

Solución: Este problema difiere del ejemplo 6.11 en que ahora asignamos una medida de 43 ohms a todos los resistores cuyas resistencias sean mayores que 42.5 y menores que 43.5. Lo que estamos haciendo realmente es aproximar una distribución discreta por medio de una distribución continua normal. El área que se requiere es la región sombreada a la derecha de 43.5 en la figura 6.19. Encontramos ahora que

$$z = \frac{43.5 - 40}{2} = 1.75.$$

En consecuencia,

$$P(X > 43.5) = P(Z > 1.75) = 1 - P(Z < 1.75) = 1 - 0.9599 = 0.0401.$$

Por lo tanto, 4.01% de las resistencias exceden 43 ohms cuando se miden al ohm más cercano. La diferencia $6.68\% - 4.01\% = 2.67\%$ entre esta respuesta y la del ejemplo 6.11 representa todos los valores de resistencias mayores que 43 y menores que 43.5, que ahora se registran como de 43 ohms. ┘

Ejemplo 6.13: La calificación promedio para un examen es 74 y la desviación estándar es 7. Si 12% del grupo obtiene A y las calificaciones siguen una curva que tiene una distribución normal, ¿cuál es la A más baja posible y la B más alta posible?

Solución: En este ejemplo comenzamos con un área de probabilidad conocida, calculamos el valor z y después determinamos x con la fórmula $x = \sigma z + \mu$. Un área de 0.12, que corresponde a la fracción de estudiantes que reciben A , está sombreada en la figura 6.20. Necesitamos un valor z que deje 0.12 del área a la derecha y, por lo tanto, un área de 0.88 a la izquierda. A partir de la tabla A.3, $P(Z < 1.18)$ tiene el valor más cercano a 0.88, de manera que el valor z que se desea es 1.18. En consecuencia,

$$x = (7)(1.18) + 74 = 82.26.$$

Por lo tanto, la A más baja es 83 y la B más alta es 82. ▮

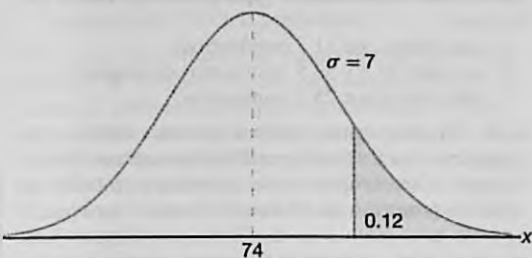


Figura 6.20: Área para el ejemplo 6.13.

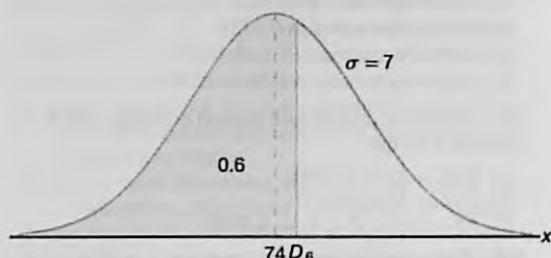


Figura 6.21: Área para el ejemplo 6.14.

Ejemplo 6.14: Remítase al ejemplo 6.13 y calcule el sexto decil.

Solución: El sexto decil, escrito como D_6 , es el valor x que deja 60% del área a la izquierda, como se muestra en la figura 6.21. En la tabla A.3 encontramos que $P(Z < 0.25) \approx 0.6$, de manera que el valor z deseado es 0.25. Ahora, $x = (7)(0.25) + 74 = 75.75$. Por lo tanto, $D_6 = 75.75$. Es decir, 60% de las calificaciones son 75 o menos. ▮

Ejercicios

6.1 Dada una distribución continua uniforme, demuestre que

a) $\mu = \frac{A+B}{2}$, y

b) $\sigma^2 = \frac{(B-A)^2}{12}$.

6.2 Suponga que X tiene una distribución continua uniforme de 1 a 5. Determine la probabilidad condicional $P(X > 2.5 \mid X \leq 4)$.

6.3 La cantidad de café diaria, en litros, que sirve una máquina que se localiza en el vestíbulo de un aeropuerto es una variable aleatoria X que tiene una

distribución continua uniforme con $A = 7$ y $B = 10$. Calcule la probabilidad de que en un día determinada la cantidad de café que sirve esta máquina sea

a) a lo sumo 8.8 litros;

b) más de 7.4 litros, pero menos de 9.5 litros;

c) al menos 8.5 litros.

6.4 Un autobús llega cada 10 minutos a una parada. Se supone que el tiempo de espera para un individuo en particular es una variable aleatoria con distribución continua uniforme.

- a) ¿Cuál es la probabilidad de que el individuo espere más de 7 minutos?
 b) ¿Cuál es la probabilidad de que el individuo espere entre 2 y 7 minutos?

6.5 Dada una distribución normal estándar, calcule el área bajo la curva que está

- a) a la izquierda de $z = -1.39$;
 b) a la derecha de $z = 1.96$;
 c) entre $z = -2.16$ y $z = -0.65$;
 d) a la izquierda de $z = 1.43$;
 e) a la derecha de $z = -0.89$;
 f) entre $z = -0.48$ y $z = 1.74$.

6.6 Calcule el valor de z si el área bajo una curva normal estándar

- a) a la derecha de z es 0.3622;
 b) a la izquierda de z es 0.1131;
 c) entre 0 y z , con $z > 0$, es 0.4838;
 d) entre $-z$ y z , con $z > 0$, es 0.9500.

6.7 Dada una distribución normal estándar, calcule el valor de k tal que

- a) $P(Z > k) = 0.2946$;
 b) $P(Z < k) = 0.0427$;
 c) $P(-0.93 < Z < k) = 0.7235$.

6.8 Dada una distribución normal con $\mu = 30$ y $\sigma = 6$, calcule

- a) el área de la curva normal a la derecha de $x = 17$;
 b) el área de la curva normal a la izquierda de $x = 22$;
 c) el área de la curva normal entre $x = 32$ y $x = 41$;
 d) el valor de x que tiene 80% del área de la curva normal a la izquierda;
 e) los dos valores de x que contienen 75% central del área de la curva normal.

6.9 Dada la variable X normalmente distribuida con una media de 18 y una desviación estándar de 2.5, calcule

- a) $P(X < 15)$;
 b) el valor de k tal que $P(X < k) = 0.2236$;
 c) el valor de k tal que $P(X > k) = 0.1814$;
 d) $P(17 < X < 21)$.

6.10 De acuerdo con el teorema de Chebyshev, la probabilidad de que cualquier variable aleatoria tome un valor dentro de 3 desviaciones estándar de la media es de al menos $8/9$. Si se sabe que la distribución de probabilidad de una variable aleatoria X es normal con media μ y varianza σ^2 , ¿cuál es el valor exacto de $P(\mu - 3\sigma < X < \mu + 3\sigma)$?

6.11 Una máquina expendedora de bebidas gaseosas se regula para que sirva un promedio de 200 mililitros por vaso. Si la cantidad de bebida se distribuye nor-

malmente con una desviación estándar igual a 15 mililitros,

- a) ¿qué fracción de los vasos contendrá más de 224 mililitros?
 b) ¿cuál es la probabilidad de que un vaso contenga entre 191 y 209 mililitros?
 c) ¿cuántos vasos probablemente se derramarán si se utilizan vasos de 230 mililitros para las siguientes 1000 bebidas?
 d) ¿por debajo de qué valor obtendremos el 25% más bajo en el llenado de las bebidas?

6.12 Las barras de pan de centeno que cierta panadería distribuye a las tiendas locales tienen una longitud promedio de 30 centímetros y una desviación estándar de 2 centímetros. Si se supone que las longitudes están distribuidas normalmente, ¿qué porcentaje de las barras son

- a) más largas que 31.7 centímetros?
 b) de entre 29.3 y 33.5 centímetros de longitud?
 c) más cortas que 25.5 centímetros?

6.13 Un investigador informa que unos ratones a los que primero se les restringen drásticamente sus dietas y después se les enriquecen con vitaminas y proteínas vivirán un promedio de 40 meses. Si suponemos que la vida de tales ratones se distribuye normalmente, con una desviación estándar de 6.3 meses, calcule la probabilidad de que un ratón determinado viva

- a) más de 32 meses;
 b) menos de 28 meses;
 c) entre 37 y 49 meses.

6.14 El diámetro interior del anillo de un pistón terminado se distribuye normalmente con una media de 10 centímetros y una desviación estándar de 0.03 centímetros.

- a) ¿Qué proporción de anillos tendrá diámetros interiores que excedan 10.075 centímetros?
 b) ¿Cuál es la probabilidad de que el anillo de un pistón tenga un diámetro interior de entre 9.97 y 10.03 centímetros?
 c) ¿Por debajo de qué valor del diámetro interior caerá el 15% de los anillos de pistón?

6.15 Un abogado viaja todos los días de su casa en los suburbios a su oficina en el centro de la ciudad. El tiempo promedio para un viaje sólo de ida es de 24 minutos, con una desviación estándar de 3.8 minutos. Si se supone que la distribución de los tiempos de viaje está distribuida normalmente.

- a) ¿Cuál es la probabilidad de que un viaje tome al menos 1/2 hora?
 b) Si la oficina abre a las 9:00 A.M. y él sale diario de su casa a las 8:45 A.M., ¿qué porcentaje de las veces llegará tarde al trabajo?

- c) Si sale de su casa a las 8:35 A.M. y el café se sirve en la oficina de 8:50 A.M. a 9:00 A.M., ¿cuál es la probabilidad de que se pierda el café?
- d) Calcule la duración mayor en la que se encuentra el 15% de los viajes más lentos.
- e) Calcule la probabilidad de que 2 de los siguientes 3 viajes tomen al menos 1/2 hora.

6.16 En el ejemplar de noviembre de 1990 de *Chemical Engineering Progress*, un estudio analiza el porcentaje de pureza del oxígeno de cierto proveedor. Suponga que la media fue de 99.61, con una desviación estándar de 0.08. Suponga que la distribución del porcentaje de pureza fue aproximadamente normal.

- a) ¿Qué porcentaje de los valores de pureza esperaría que estuvieran entre 99.5 y 99.7?
- b) ¿Qué valor de pureza esperaría que excediera exactamente 5% de la población?

6.17 La vida promedio de cierto tipo de motor pequeño es de 10 años, con una desviación estándar de 2 años. El fabricante reemplaza gratis todos los motores que fallen dentro del periodo de garantía. Si estuviera dispuesto a reemplazar sólo 3% de los motores que fallan, ¿cuánto tiempo de garantía debería ofrecer? Suponga que la duración de un motor sigue una distribución normal.

6.18 La estatura de 1000 estudiantes se distribuye normalmente con una media de 174.5 centímetros y una desviación estándar de 6.9 centímetros. Si se supone que las estaturas se redondean al medio centímetro más cercano, ¿cuántos de estos estudiantes esperaría que tuvieran una estatura

- a) menor que 160.0 centímetros?
- b) de entre 171.5 y 182.0 centímetros inclusive?
- c) igual a 175.0 centímetros?
- d) mayor o igual que 188.0 centímetros?

6.19 Una empresa paga a sus empleados un salario promedio de \$15.90 por hora, con una desviación estándar de \$1.50. Si los salarios se distribuyen aproximadamente de forma normal y se redondean al centavo más cercano,

- a) ¿qué porcentaje de los trabajadores recibe salarios de entre \$13.75 y \$16.22 por hora?
- b) ¿el 5% de los salarios más altos por hora de los empleados es mayor a qué cantidad?

6.20 Los pesos de un gran número de *poodle* miniatura se distribuyen aproximadamente de forma normal con una media de 8 kilogramos y una desviación estándar de 0.9 kilogramos. Si las mediciones se redondean al décimo de kilogramo más cercano, calcule la fracción de estos *poodle* con pesos

- a) por arriba de 9.5 kilogramos;
- b) a lo sumo 8.6 kilogramos;
- c) entre 7.3 y 9.1 kilogramos.

6.21 La resistencia a la tensión de cierto componente de metal se distribuye normalmente con una media de 10,000 kilogramos por centímetro cuadrado y una desviación estándar de 100 kilogramos por centímetro cuadrado. Las mediciones se redondean a los 50 kilogramos por centímetro cuadrado más cercanos.

- a) ¿Qué proporción de estos componentes excede a 10,150 kilogramos por centímetro cuadrado de resistencia a la tensión?
- b) Si las especificaciones requieren que todos los componentes tengan una resistencia a la tensión de entre 9800 y 10,200 kilogramos por centímetro cuadrado, ¿qué proporción de piezas esperaría que se descartara?

6.22 Si un conjunto de observaciones se distribuye de manera normal, ¿qué porcentaje de éstas difieren de la media en

- a) más de 1.3σ ?
- b) menos de 0.52σ ?

6.23 El coeficiente intelectual (CI) de 600 aspirantes a cierta universidad se distribuye aproximadamente de forma normal con una media de 115 y una desviación estándar de 12. Si la universidad requiere un CI de al menos 95, ¿cuántos de estos estudiantes serán rechazados con base en éste sin importar sus otras calificaciones? Tome en cuenta que el CI de los aspirantes se redondea al entero más cercano.

6.5 Aproximación normal a la binomial

Las probabilidades asociadas con experimentos binomiales se obtienen fácilmente a partir de la fórmula $b(x; n, p)$ de la distribución binomial o de la tabla A.1 cuando n es pequeña. Además, las probabilidades binomiales están disponibles en muchos paquetes de software. Sin embargo, resulta aleccionador conocer la relación entre la distribución binomial y la normal. En la sección 5.5 explicamos cómo se puede utilizar la distribución de Poisson para aproximar probabilidades binomiales cuando n es muy grande y p se acerca mucho a 0 o a 1. Tanto la distribución binomial como la de Poisson son

discretas. La primera aplicación de una distribución continua de probabilidad para aproximar probabilidades sobre un espacio muestral discreto se demostró en el ejemplo 6.12, donde se utilizó la curva normal. La distribución normal a menudo es una buena aproximación a una distribución discreta cuando la última adquiere una forma de campana simétrica. Desde un punto de vista teórico, algunas distribuciones convergen a la normal a medida que sus parámetros se aproximan a ciertos límites. La distribución normal es una distribución de aproximación conveniente, ya que la función de distribución acumulativa se tabula con mucha facilidad. La distribución binomial se aproxima bien por medio de la normal en problemas prácticos cuando se trabaja con la función de distribución acumulativa. Ahora plantearemos un teorema que nos permitirá utilizar áreas bajo la curva normal para aproximar propiedades binomiales cuando n es suficientemente grande.

Teorema 6.3: Si X es una variable aleatoria binomial con media $\mu = np$ y varianza $\sigma^2 = npq$, entonces la forma limitante de la distribución de

$$Z = \frac{X - np}{\sqrt{npq}},$$

conforme $n \rightarrow \infty$, es la distribución normal estándar $n(z; 0, 1)$.

Resulta que la distribución normal con $\mu = np$ y $\sigma^2 = np(1 - p)$ no sólo ofrece una aproximación muy precisa a la distribución binomial cuando n es grande y p no está extremadamente cerca de 0 o de 1, sino que también brinda una aproximación bastante buena aun cuando n es pequeña y p está razonablemente cerca de $1/2$.

Para ilustrar la aproximación normal a la distribución binomial primero dibujamos el histograma para $b(x; 15, 0.4)$ y después superponemos la curva normal particular con la misma media y varianza que la variable binomial X . En consecuencia, dibujamos una curva normal con

$$\mu = np = (15)(0.4) = 6 \text{ y } \sigma^2 = npq = (15)(0.4)(0.6) = 3.6.$$

El histograma de $b(x; 15, 0.4)$ y la curva normal superpuesta correspondiente, que está determinada por completo por su media y su varianza, se ilustran en la figura 6.22.

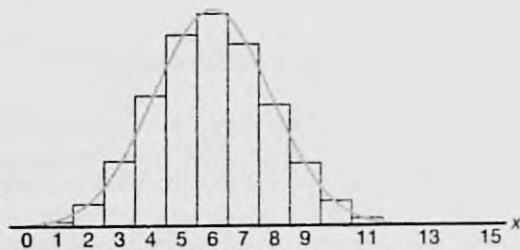


Figura 6.22: Aproximación normal de $b(x; 15, 0.4)$.

La probabilidad exacta de que la variable aleatoria binomial X tome un valor determinado x es igual al área de la barra cuya base se centra en x . Por ejemplo, la probabilidad exacta de que X tome el valor 4 es igual al área del rectángulo con base centrada en $x = 4$. Si usamos la tabla A.1, encontramos que esta área es

$$P(X = 4) = b(4; 15, 0.4) = 0.1268,$$

que es aproximadamente igual al área de la región sombreada bajo la curva normal entre las dos ordenadas $x_1 = 3.5$ y $x_2 = 4.5$ en la figura 6.23. Al convertir a valores z , tenemos

$$z_1 = \frac{3.5 - 6}{1.897} = -1.32 \quad \text{y} \quad z_2 = \frac{4.5 - 6}{1.897} = -0.79.$$

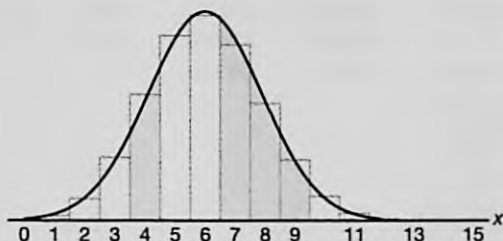


Figura 6.23: Aproximación normal de $b(x; 15, 0.4)$ y $\sum_{x=7}^9 b(x; 15, 0.4)$.

Si X es una variable aleatoria binomial y Z una variable normal estándar, entonces,

$$\begin{aligned} P(X = 4) &= b(4; 15, 0.4) \approx P(-1.32 < Z < -0.79) \\ &= P(Z < -0.79) - P(Z < -1.32) = 0.2148 - 0.0934 = 0.1214. \end{aligned}$$

Esto se aproxima bastante al valor exacto de 0.1268.

La aproximación normal es más útil en el cálculo de sumatorias binomiales para valores grandes de n . Si nos remitimos a la figura 6.23, nos podríamos interesar en la probabilidad de que X tome un valor de 7 a 9. La probabilidad exacta es dada por

$$\begin{aligned} P(7 \leq X \leq 9) &= \sum_{x=7}^9 b(x; 15, 0.4) - \sum_{x=0}^6 b(x; 15, 0.4) \\ &= 0.9662 - 0.6098 = 0.3564, \end{aligned}$$

que es igual a la sumatoria de las áreas de los rectángulos cuyas bases están centradas en $x = 7, 8$ y 9 . Para la aproximación normal calculamos el área de la región sombreada bajo la curva entre las ordenadas $x_1 = 6.5$ y $x_2 = 9.5$ de la figura 6.23. Los valores z correspondientes son

$$z_1 = \frac{6.5 - 6}{1.897} = 0.26 \quad \text{y} \quad z_2 = \frac{9.5 - 6}{1.897} = 1.85.$$

Ahora,

$$\begin{aligned} P(7 \leq X \leq 9) &\approx P(0.26 < Z < 1.85) = P(Z < 1.85) - P(Z < 0.26) \\ &= 0.9678 - 0.6026 = 0.3652. \end{aligned}$$

Una vez más, la aproximación de la curva normal ofrece un valor que se acerca al valor exacto de 0.3564. El grado de exactitud, que depende de qué tan bien se ajuste la curva al histograma, se incrementa a medida que aumenta n . Esto es particularmente cierto cuando p no está muy cerca de $1/2$ y el histograma ya no es simétrico. Las figuras 6.24 y 6.25 muestran los histogramas para $b(x; 6, 0.2)$ y $b(x; 15, 0.2)$, respectivamente. Es evidente que una curva normal se ajustará mucho mejor al histograma cuando $n = 15$ que cuando $n = 6$.

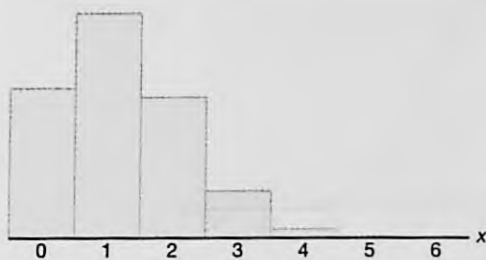


Figura 6.24: Histograma para $b(x; 6, 0.2)$.

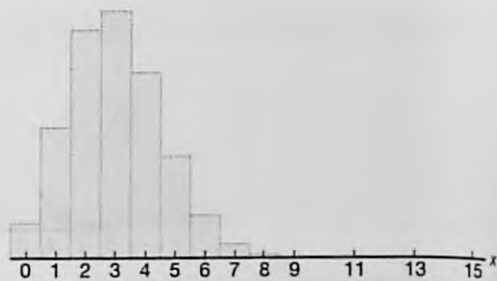


Figura 6.25: Histograma para $b(x; 15, 0.2)$.

En las ilustraciones de la aproximación normal a la binomial se hizo evidente que si buscamos el área bajo la curva normal hacia la izquierda de, digamos x , es más preciso utilizar $x + 0.5$. Esto es una corrección para dar cabida al hecho de que una distribución discreta se aproxima mediante una distribución continua. La corrección $+0.5$ se llama **corrección de continuidad**. La explicación anterior conduce a la siguiente aproximación normal formal a la binomial.

Aproximación normal a la distribución binomial Sea X una variable aleatoria binomial con parámetros n y p . Para una n grande, X tiene aproximadamente una distribución normal con $\mu = np$ y $\sigma^2 = npq = np(1-p)$ y

$$\begin{aligned} P(X \leq x) &= \sum_{k=0}^x b(k; n, p) \\ &\approx \text{área bajo la curva normal a la izquierda de } x + 0.5 \\ &= P\left(Z \leq \frac{x + 0.5 - np}{\sqrt{npq}}\right), \end{aligned}$$

y la aproximación será buena si np y $n(1-p)$ son mayores que o iguales a 5.

Como indicamos antes, la calidad de la aproximación es muy buena para n grande. Si p está cerca de $1/2$, un tamaño de la muestra moderado o pequeño será suficiente para una aproximación razonable. Ofrecemos la tabla 6.1 como una indicación de la calidad

de la aproximación. Se presentan tanto la aproximación normal como las probabilidades binomiales acumulativas reales. Observe que en $p = 0.05$ y $p = 0.10$ la aproximación es muy burda para $n = 10$. Sin embargo, incluso para $n = 10$, observe la mejoría para $p = 0.50$. Por otro lado, cuando p es fija en $p = 0.05$, observe cómo mejora la aproximación conforme vamos de $n = 20$ a $n = 100$.

Tabla 6.1: Aproximación normal y probabilidades binomiales acumulativas reales

r	$p = 0.05, n = 10$		$p = 0.10, n = 10$		$p = 0.50, n = 10$	
	Binomial	Normal	Binomial	Normal	Binomial	Normal
0	0.5987	0.5000	0.3487	0.2981	0.0010	0.0022
1	0.9139	0.9265	0.7361	0.7019	0.0107	0.0136
2	0.9885	0.9981	0.9298	0.9429	0.0547	0.0571
3	0.9990	1.0000	0.9872	0.9959	0.1719	0.1711
4	1.0000	1.0000	0.9984	0.9999	0.3770	0.3745
5			1.0000	1.0000	0.6230	0.6255
6					0.8281	0.8289
7					0.9453	0.9429
8					0.9893	0.9864
9					0.9990	0.9978
10					1.0000	0.9997

r	$p = 0.05$					
	$n = 20$		$n = 50$		$n = 100$	
	Binomial	Normal	Binomial	Normal	Binomial	Normal
0	0.3585	0.3015	0.0769	0.0968	0.0059	0.0197
1	0.7358	0.6985	0.2794	0.2578	0.0371	0.0537
2	0.9245	0.9382	0.5405	0.5000	0.1183	0.1251
3	0.9841	0.9948	0.7604	0.7422	0.2578	0.2451
4	0.9974	0.9998	0.8964	0.9032	0.4360	0.4090
5	0.9997	1.0000	0.9622	0.9744	0.6160	0.5910
6	1.0000	1.0000	0.9882	0.9953	0.7660	0.7549
7			0.9968	0.9994	0.8720	0.8749
8			0.9992	0.9999	0.9369	0.9463
9			0.9998	1.0000	0.9718	0.9803
10			1.0000	1.0000	0.9885	0.9941

Ejemplo 6.15: Un paciente que padece una rara enfermedad de la sangre tiene 0.4 de probabilidad de recuperarse. Si se sabe que 100 personas contrajeron esta enfermedad, ¿cuál es la probabilidad de que sobrevivan menos de 30?

Solución: Representemos con la variable binomial X el número de pacientes que sobreviven. Como $n = 100$, deberíamos obtener resultados muy precisos usando la aproximación de la curva normal con

$$\mu = np = (100)(0.4) = 40 \quad \text{y} \quad \sigma = \sqrt{npq} = \sqrt{(100)(0.4)(0.6)} = 4.899.$$

Para obtener la probabilidad que se desea, tenemos que calcular el área a la izquierda de $x = 29.5$.

El valor z que corresponde a 29.5 es

$$z = \frac{29.5 - 40}{4.899} = -2.14,$$

y la probabilidad de que menos de 30 de los 100 pacientes sobrevivan está dada por la región sombreada en la figura 6.26. Por lo tanto,

$$P(X < 30) \approx P(Z < -2.14) = 0.0162.$$

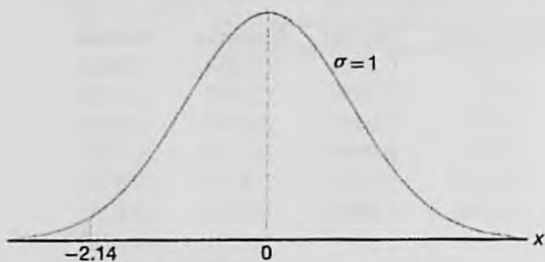


Figura 6.26: Área para el ejemplo 6.15.

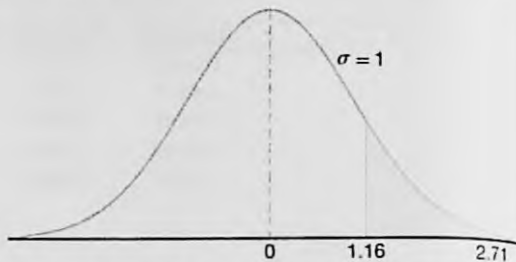


Figura 6.27: Área para el ejemplo 6.16.

Ejemplo 6.16: Un examen de opción múltiple tiene 200 preguntas, cada una con 4 respuestas posibles, de las que sólo una es la correcta. ¿Cuál es la probabilidad de que solamente adivinando se obtengan de 25 a 30 respuestas correctas para 80 de los 200 problemas sobre los que el estudiante no tiene conocimientos?

Solución: La probabilidad de adivinar una respuesta correcta para cada una de las 80 preguntas es $p = 1/4$. Si X representa el número de respuestas correctas sólo porque se adivinaron, entonces,

$$P(25 \leq X \leq 30) = \sum_{x=25}^{30} b(x; 80, 1/4).$$

Al usar la aproximación de la curva normal con

$$\mu = np = (80) \left(\frac{1}{4} \right) = 20$$

y

$$\sigma = \sqrt{npq} = \sqrt{(80)(1/4)(3/4)} = 3.873,$$

necesitamos el área entre $x_1 = 24.5$ y $x_2 = 30.5$. Los valores z correspondientes son

$$z_1 = \frac{24.5 - 20}{3.873} = 1.16 \text{ y } z_2 = \frac{30.5 - 20}{3.873} = 2.71.$$

La probabilidad de adivinar correctamente de 25 a 30 preguntas es dada por la región sombreada de la figura 6.27. En la tabla A.3 encontramos que

$$\begin{aligned} P(25 \leq X \leq 30) &= \sum_{x=25}^{30} b(x; 80, 0.25) \approx P(1.16 < Z < 2.71) \\ &= P(Z < 2.71) - P(Z < 1.16) = 0.9966 - 0.8770 = 0.1196. \end{aligned}$$

Ejercicios

6.24 Se lanza una moneda 400 veces. Utilice la aproximación a la curva normal para calcular la probabilidad de obtener

- entre 185 y 210 caras;
- exactamente 205 caras;
- menos de 176 o más de 227 caras.

6.25 En un proceso para fabricar un componente electrónico, 1% de los artículos resultan defectuosos. Un plan de control de calidad consiste en seleccionar 100 artículos de un proceso de producción y detenerlo o continuar con él si ninguno está defectuoso. Use la aproximación normal a la binomial para calcular

- la probabilidad de que el proceso continúe con el plan de muestreo descrito;
- la probabilidad de que el proceso continúe aun si éste va mal (es decir, si la frecuencia de componentes defectuosos cambió a 5.0% de defectuosos).

6.26 Un proceso produce 10% de artículos defectuosos. Si se seleccionan al azar 100 artículos del proceso, ¿cuál es la probabilidad de que el número de defectuosos

- exceda los 13?
- sea menor que 8?

6.27 Un paciente tiene 0.9 de probabilidad de recuperarse de una operación de corazón delicada. De los siguientes 100 pacientes que se someten a esta operación, ¿cuál es la probabilidad de que

- sobrevivan entre 84 y 95 inclusive?
- sobrevivan menos de 86?

6.28 Investigadores de la Universidad George Washington y del Instituto Nacional de Salud informan que aproximadamente 75% de las personas cree que "los tranquilizantes funcionan muy bien para lograr que una persona esté más tranquila y relajada". De las siguientes 80 personas entrevistadas, ¿cuál es la probabilidad de que

- al menos 50 tengan esta opinión?
- a lo sumo 56 tengan esta opinión?

6.29 Si 20% de los residentes de una ciudad de Estados Unidos prefieren un teléfono blanco sobre cualquier otro color disponible, ¿cuál es la probabilidad de que, de los siguientes 1000 teléfonos que se instalen en esa ciudad,

- entre 170 y 185 sean blancos?
- al menos 210 pero no más de 225 sean blancos?

6.30 Un fabricante de medicamentos sostiene que cierto medicamento cura una enfermedad de la sangre, en promedio, 80% de las veces. Para verificar la aseveración, inspectores gubernamentales utilizan el medi-

camento en una muestra de 100 individuos y deciden aceptar la afirmación si se curan 75 o más.

- ¿Cuál es la probabilidad de que los inspectores gubernamentales rechacen la aseveración si la probabilidad de curación es, de hecho, de 0.8?
- ¿Cuál es la probabilidad de que el gobierno acepte la afirmación si la probabilidad de curación resulta tan baja como 0.7?

6.31 Una sexta parte de los estudiantes de primer año que entran a una escuela estatal grande provienen de otros estados. Si son asignados al azar a los 180 dormitorios de un edificio, ¿cuál es la probabilidad de que en un determinado dormitorio al menos una quinta parte de los estudiantes provenga de otro estado?

6.32 Una empresa farmacéutica sabe que aproximadamente 5% de sus píldoras anticonceptivas no contiene la cantidad suficiente de un ingrediente, lo que las vuelve ineficaces. ¿Cuál es la probabilidad de que menos de 10 píldoras en una muestra de 200 sean ineficaces?

6.33 Estadísticas publicadas por la National Highway Traffic Safety Administration y el National Safety Council revelan que en una noche promedio de fin de semana, uno de cada 10 conductores está ebrio. Si la siguiente noche de sábado se revisan 400 conductores al azar, ¿cuál es la probabilidad de que el número de conductores ebrios sea

- menor que 32?
- mayor que 49?
- al menos 35 pero menos que 47?

6.34 Un par de dados se lanza 180 veces. ¿Cuál es la probabilidad de que ocurra un total de 7

- al menos 25 veces?
- entre 33 y 41 veces?
- exactamente 30 veces?

6.35 Una empresa produce partes componentes para un motor. Las especificaciones de las partes sugieren que sólo 95% de los artículos las cumplen. Las partes para los clientes se embarcan en lotes de 100.

- ¿Cuál es la probabilidad de que más de 2 artículos estén defectuosos en un lote determinado?
- ¿Cuál es la probabilidad de que más de 10 artículos de un lote estén defectuosos?

6.36 Una práctica común por parte de las aerolíneas consiste en vender más boletos que el número real de asientos para un vuelo específico porque los clientes que compran boletos no siempre se presentan a abordar el avión. Suponga que el porcentaje de pasajeros que no se presentan a la hora del vuelo es de 2%. Para un vuelo particular con 197 asientos, se vendieron un total

de 200 boletos. ¿Cuál es la probabilidad de que la aerolínea haya sobrevendido el vuelo?

6.37 El nivel X de colesterol en la sangre en muchachos de 14 años tiene aproximadamente una distribución normal, con una media de 170 y una desviación estándar de 30.

- Determine la probabilidad de que el nivel de colesterol en la sangre de un muchacho de 14 años elegido al azar exceda 230.
- En una escuela secundaria hay 300 muchachos de 14 años. Determine la probabilidad de que por lo menos 8 de ellos tengan un nivel de colesterol superior a 230.

6.38 Una empresa de telemarketing tiene una máquina especial para abrir cartas que abre y extrae el contenido de los sobres. Si un sobre se colocara de forma incorrecta en la máquina, no se podría extraer su contenido, o incluso se podría dañar. En este caso se dice que "falló" la máquina.

- Si la probabilidad de que falle la máquina es de 0.01, ¿cuál es la probabilidad de que ocurra más de una falla en un lote de 20 sobres?
- Si la probabilidad de que falle la máquina es de 0.01 y se abrirá un lote de 500 sobres, ¿cuál es la probabilidad de que ocurran más de 8 fallas?

6.6 Distribución gamma y distribución exponencial

Aunque la distribución normal se puede utilizar para resolver muchos problemas de ingeniería y ciencias, aún hay numerosas situaciones que requieren diferentes tipos de funciones de densidad. En esta sección se estudiarán dos de estas funciones de densidad, la **distribución gamma** y la **distribución exponencial**.

Resulta que la distribución exponencial es un caso especial de la distribución gamma, y ambas tienen un gran número de aplicaciones. La distribución exponencial y la distribución gamma desempeñan un papel importante en la teoría de colas y en problemas de confiabilidad. Los tiempos entre llegadas en instalaciones de servicio y los tiempos de operación antes de que partes componentes y sistemas eléctricos empiecen a fallar a menudo se representan bien mediante la distribución exponencial. La relación entre la distribución gamma y la exponencial permite que la gamma se utilice en problemas similares. En la siguiente sección se presentarán más detalles y ejemplos.

La distribución gamma deriva su nombre de la bien conocida **función gamma**, que se estudia en muchas áreas de las matemáticas. Antes de estudiar la distribución gamma repasaremos esta función y algunas de sus propiedades importantes.

Definición 6.2: La función gamma se define como

$$\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx, \text{ para } \alpha > 0.$$

Las siguientes son algunas propiedades sencillas de la función gamma.

- $\Gamma(n) = (n-1)(n-2) \cdots (1) \Gamma(1)$ para una integral positiva n .

Para ver la demostración, al integrar por partes con $u = x^{\alpha-1}$ y $dv = e^{-x} dx$, obtenemos

$$\Gamma(\alpha) = -e^{-x} x^{\alpha-1} \Big|_0^{\infty} + \int_0^{\infty} e^{-x} (\alpha-1)x^{\alpha-2} dx = (\alpha-1) \int_0^{\infty} x^{\alpha-2} e^{-x} dx,$$

para $\alpha > 1$, que produce la fórmula recursiva

$$\Gamma(\alpha) = (\alpha-1)\Gamma(\alpha-1).$$

El resultado proviene de la aplicación repetida de la fórmula recursiva. Si utilizamos este resultado, podemos demostrar con facilidad las siguientes dos propiedades.

b) $\Gamma(n) = (n - 1)!$ para una integral positiva n .

c) $\Gamma(1) = 1$.

Asimismo, tenemos la siguiente propiedad de $\Gamma(\alpha)$, que el lector deberá verificar (véase el ejercicio 6.39 de la página 206).

d) $\Gamma(1/2) = \sqrt{\pi}$.

A continuación se define la **distribución gamma**.

Distribución gamma La variable aleatoria continua X tiene una **distribución gamma**, con parámetros α y β , si su función de densidad está dada por

$$f(x; \alpha, \beta) = \begin{cases} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta}, & x > 0, \\ 0, & \text{en otro caso,} \end{cases}$$

donde $\alpha > 0$ y $\beta > 0$.

En la figura 6.28 se muestran gráficas de varias distribuciones gamma para ciertos valores específicos de los parámetros α y β . La distribución gamma especial para la que $\alpha = 1$ se llama **distribución exponencial**.

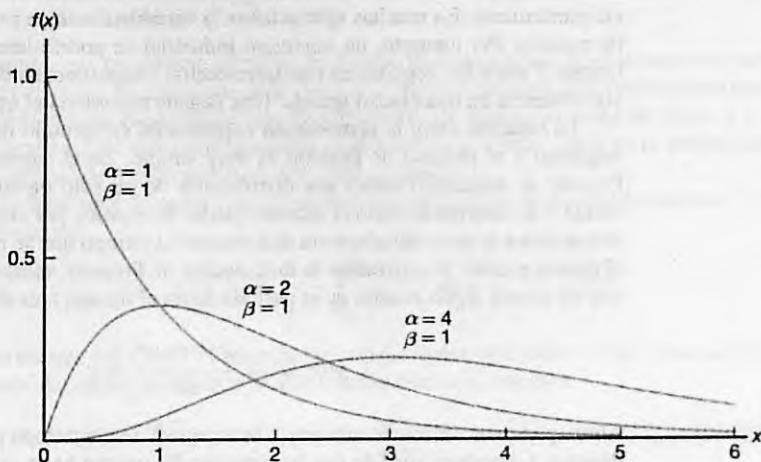


Figura 6.28: Distribuciones gamma.

Distribución exponencial La variable aleatoria continua X tiene una **distribución exponencial**, con parámetro β , si su función de densidad es dada por

$$f(x; \beta) = \begin{cases} \frac{1}{\beta} e^{-x/\beta}, & x > 0, \\ 0, & \text{en otro caso,} \end{cases}$$

donde $\beta > 0$.

El siguiente teorema y corolario proporcionan la media y la varianza de la distribución gamma y la exponencial.

Teorema 6.4: La media y la varianza de la distribución gamma son

$$\mu = \alpha\beta \text{ y } \sigma^2 = \alpha\beta^2.$$

La demostración de este teorema se encuentra en el apéndice A.26.

Corolario 6.1: La media y la varianza de la distribución exponencial son

$$\mu = \beta \text{ y } \sigma^2 = \beta^2.$$

Relación con el proceso de Poisson

Continuaremos con las aplicaciones de la distribución exponencial y después regresaremos a la distribución gamma. Las aplicaciones más importantes de la distribución exponencial son situaciones donde se aplica el proceso de Poisson (véase la sección 5.5). El lector debería recordar que el proceso de Poisson permite utilizar la distribución discreta llamada distribución de Poisson. Recuerde que la distribución de Poisson se utiliza para calcular la probabilidad de números específicos de “eventos” durante un *periodo o espacio* particulares. En muchas aplicaciones la variable aleatoria es el tiempo o la cantidad de espacio. Por ejemplo, un ingeniero industrial se podría interesar en un modelo de tiempo T entre las llegadas en una intersección congestionada durante las horas de mayor afluencia en una ciudad grande. Una llegada representa el evento de Poisson.

La relación entre la distribución exponencial (a menudo denominada exponencial negativa) y el proceso de Poisson es muy simple. En el capítulo 5 la distribución de Poisson se desarrolló como una distribución de un solo parámetro con parámetro λ , donde λ se interpreta como el número medio de eventos por *unidad de “tiempo”*. Considere ahora la *variable aleatoria* descrita por el tiempo que se requiere para que ocurra el primer evento. Si utilizamos la distribución de Poisson, vemos que la probabilidad de que no ocurra algún evento, en el periodo hasta el tiempo t , es dada por

$$p(0; \lambda t) = \frac{e^{-\lambda t} (\lambda t)^0}{0!} = e^{-\lambda t}.$$

Ahora podemos utilizar lo anterior y hacer que X sea el tiempo para el primer evento de Poisson. La probabilidad de que la duración del tiempo hasta el primer evento exceda x es la misma que la probabilidad de que no ocurra algún evento de Poisson en x . Esto último, por supuesto, es dado por $e^{-\lambda x}$. Como resultado,

$$P(X > x) = e^{-\lambda x}.$$

Así, la función de distribución acumulativa para X es dada por

$$P(0 \leq X \leq x) = 1 - e^{-\lambda x}.$$

Ahora, para poder reconocer la presencia de la distribución exponencial, podemos diferenciar la función de distribución acumulativa anterior con el fin de obtener la función de densidad

$$f(x) = \lambda e^{-\lambda x},$$

que es la función de densidad de la distribución exponencial con $\lambda = 1/\beta$.

Aplicaciones de la distribución exponencial y la distribución gamma

En la explicación anterior establecimos las bases para la aplicación de la distribución exponencial en el “tiempo de llegada” o tiempo para problemas con eventos de Poisson. Aquí ilustraremos algunas aplicaciones de modelado y después procederemos a analizar el papel que la distribución gamma desempeña en ellas. Observe que la media de la distribución exponencial es el parámetro β , el recíproco del parámetro en la distribución de Poisson. El lector debería recordar que con frecuencia se dice que la distribución de Poisson no tiene memoria, lo cual implica que las ocurrencias en periodos sucesivos son independientes. El importante parámetro β es el tiempo promedio entre eventos. En la teoría de confiabilidad, donde la falla de equipo con frecuencia se ajusta a este proceso de Poisson, β se denomina **tiempo medio entre fallas**. Muchas descomposturas de equipo siguen el proceso de Poisson y por ello se aplica la distribución exponencial. Otras aplicaciones incluyen tiempos de supervivencia en experimentos biomédicos y tiempo de respuesta de computadoras.

En el siguiente ejemplo mostramos una aplicación simple de la distribución exponencial a un problema de confiabilidad. La distribución binomial también desempeña un papel en la solución.

Ejemplo 6.17: Suponga que un sistema contiene cierto tipo de componente cuyo tiempo de operación antes de fallar, en años, está dado por T . La variable aleatoria T se modela bien mediante la distribución exponencial con tiempo medio de operación antes de fallar $\beta = 5$. Si se instalan 5 de estos componentes en diferentes sistemas, ¿cuál es la probabilidad de que al final de 8 años al menos dos aún funcionen?

Solución: La probabilidad de que un componente determinado siga funcionando después de 8 años es dada por

$$P(T > 8) = \frac{1}{5} \int_8^{\infty} e^{-t/5} dt = e^{-8/5} \approx 0.2.$$

Representemos con X el número de componentes que todavía funcionan después de 8 años. Entonces, utilizando la distribución binomial tenemos

$$P(X \geq 2) = \sum_{x=2}^5 b(x; 5, 0.2) = 1 - \sum_{x=0}^1 b(x; 5, 0.2) = 1 - 0.7373 = 0.2627. \quad \blacksquare$$

En el capítulo 3 se incluyen ejercicios y ejemplos en los que el lector ya se enfrentó a la distribución exponencial. Otros que implican problemas de tiempo de espera y de confiabilidad se pueden encontrar en el ejemplo 6.24 y en los ejercicios y ejercicios de repaso al final de este capítulo.

La propiedad de falta de memoria y su efecto en la distribución exponencial

En los tipos de aplicación de la distribución exponencial en los problemas de confiabilidad y de tiempo de vida de una máquina o de un componente influye la **propiedad de**

falta de memoria de la distribución exponencial. Por ejemplo, en el caso de, digamos, un componente electrónico, en el que la distribución del tiempo de vida es exponencial, la probabilidad de que el componente dure, por ejemplo, t horas, es decir, $P(X \geq t)$, es igual que la probabilidad condicional

$$P(X \geq t_0 + t \mid X \geq t_0).$$

Entonces, si el componente "alcanza" las t_0 horas, la probabilidad de que dure otras t horas es igual que la probabilidad de que dure t horas. No hay "castigo" a través del desgaste como resultado de durar las primeras t_0 horas. Por lo tanto, cuando la propiedad de falta de memoria es justificada es más adecuada la distribución exponencial. Pero si la falla del componente es resultado del desgaste lento o gradual (como en el caso del desgaste mecánico), entonces la distribución exponencial no es aplicable y serían más adecuadas la distribución gamma o la de Weibull (sección 6.10).

La importancia de la distribución gamma radica en el hecho de que define una familia en la cual otras distribuciones son casos especiales. Pero la propia distribución gamma tiene aplicaciones importantes en tiempo de espera y teoría de confiabilidad. Mientras que la distribución exponencial describe el tiempo que transcurre hasta la ocurrencia de un evento de Poisson (o el tiempo entre eventos de Poisson), el tiempo (o espacio) que transcurre hasta que *ocurre un número específico de eventos de Poisson* es una variable aleatoria, cuya función de densidad es descrita por la distribución gamma. Este número específico de eventos es el parámetro α en la función de densidad gamma. De esta manera se facilita comprender que cuando $\alpha = 1$, ocurre el caso especial de la distribución exponencial. La densidad gamma se puede desarrollar a partir de su relación con el proceso de Poisson de la misma manera en que lo hicimos con la densidad exponencial. Los detalles se dejan al lector. El siguiente es un ejemplo numérico de cómo se utiliza la distribución gamma en una aplicación de tiempo de espera.

Ejemplo 6.18: Suponga que las llamadas telefónicas que llegan a un conmutador particular siguen un proceso de Poisson con un promedio de 5 llamadas entrantes por minuto. ¿Cuál es la probabilidad de que transcurra hasta un minuto en el momento en que han entrado 2 llamadas al conmutador?

Solución: Se aplica el proceso de Poisson, con un lapso de tiempo hasta que ocurren 2 eventos de Poisson que sigue una distribución gamma con $\beta = 1/5$ y $\alpha = 2$. Denote con X el tiempo en minutos que transcurre antes de que lleguen 2 llamadas. La probabilidad que se requiere está dada por

$$P(X \leq 1) = \int_0^1 \frac{1}{\beta^2} x e^{-x/\beta} dx = 25 \int_0^1 x e^{-5x} dx = 1 - e^{-5}(1 + 5) = 0.96. \quad \blacksquare$$

Mientras el origen de la distribución gamma trata con el tiempo (o espacio) hasta la ocurrencia de α eventos de Poisson, hay muchos ejemplos donde una distribución gamma funciona muy bien aunque no exista una estructura de Poisson clara. Esto es particularmente cierto para problemas de **tiempo de supervivencia** en aplicaciones de ingeniería y biomédicas.

Ejemplo 6.19: En un estudio biomédico con ratas se utiliza una investigación de respuesta a la dosis para determinar el efecto de la dosis de un tóxico en su tiempo de supervivencia. El tóxico es producido por el combustible que utilizan los aviones y, en consecuencia, descargan con frecuencia a la atmósfera. Para cierta dosis del tóxico, el estudio determina que el tiempo de supervivencia de las ratas, en semanas, tiene una distribución gamma con $\alpha = 5$ y $\beta = 10$. ¿Cuál es la probabilidad de que una rata no sobreviva más de 60 semanas?

Solución: Sea la variable aleatoria X el tiempo de supervivencia (tiempo hasta la muerte). La probabilidad que se requiere es

$$P(X \leq 60) = \frac{1}{\beta^5} \int_0^{60} \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(5)} dx.$$

La integral anterior se puede resolver mediante la **función gamma incompleta**, que se convierte en la función de distribución acumulativa para la distribución gamma. Esta función se escribe como

$$F(x; \alpha) = \int_0^x \frac{y^{\alpha-1} e^{-y}}{\Gamma(\alpha)} dy.$$

Si permitimos que $y = x/\beta$, de modo que $x = \beta y$, tenemos

$$P(X \leq 60) = \int_0^6 \frac{y^4 e^{-y}}{\Gamma(5)} dy,$$

que se denota como $F(6; 5)$ en la tabla de la función gamma incompleta del apéndice A.23. Observe que esto permite un cálculo rápido de las probabilidades para la distribución gamma. De hecho, para este problema la probabilidad de que la rata no sobreviva más de 60 días es dada por

$$P(X \leq 60) = F(6; 5) = 0.715. \quad \blacksquare$$

Ejemplo 6.20: A partir de datos previos se sabe que la longitud de tiempo, en meses, entre las quejas de los clientes sobre cierto producto es una distribución gamma con $\alpha = 2$ y $\beta = 4$. Se realizaron cambios para hacer más estrictos los requerimientos del control de calidad después de los cuales pasaron 20 meses antes de la primera queja. ¿Parecería que los cambios realizados en el control de calidad resultaron eficaces?

Solución: Sea X el tiempo para que se presente la primera queja, el cual, en las condiciones anteriores a los cambios, seguía una distribución gamma con $\alpha = 2$ y $\beta = 4$. La pregunta se centra alrededor de qué tan raro es $X \geq 20$ dado que α y β permanecen con los valores 2 y 4, respectivamente. En otras palabras, en las condiciones anteriores ¿es razonable un "tiempo para la queja" tan grande como 20 meses? Por consiguiente, si seguimos la solución del ejemplo 6.19,

$$P(X \geq 20) = 1 - \frac{1}{\beta^\alpha} \int_0^{20} \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha)} dx.$$

De nuevo, usando $y = x/\beta$ tenemos

$$P(X \geq 20) = 1 - \int_0^5 \frac{ye^{-y}}{\Gamma(2)} dy = 1 - F(5; 2) = 1 - 0.96 = 0.04,$$

donde $F(5; 2) = 0.96$ se obtiene de la tabla A.23.

Como resultado, podríamos concluir que las condiciones de la distribución gamma con $\alpha = 2$ y $\beta = 4$ no son sustentadas por los datos de que un tiempo observado para la queja sea tan extenso como 20 meses. Entonces, es razonable concluir que el trabajo de control de calidad resultó eficaz. \blacksquare

Ejemplo 6.21: Considere el ejercicio 3.31 de la página 94. Con base en abundantes pruebas se determinó que el tiempo Y en años antes de que se requiera una reparación mayor para cierta lavadora se caracteriza por la función de densidad

$$f(y) = \begin{cases} \frac{1}{4} e^{-y/4}, & y \geq 0, \\ 0, & \text{en otro caso.} \end{cases}$$

Observe que Y es una variable aleatoria exponencial con $\mu = 4$ años. Se considera que la lavadora es una ganga si no hay probabilidades de que requiera una reparación mayor antes de cumplir 6 años de haber sido comprada. ¿Cuál es la probabilidad de $P(Y > 6)$? ¿Cuál es la probabilidad de que la lavadora requiera una reparación mayor durante el primer año?

Solución: Considere la función de distribución acumulativa $F(y)$ para la distribución exponencial,

$$F(y) = \frac{1}{\beta} \int_0^y e^{-t/\beta} dt = 1 - e^{-y/\beta}.$$

De manera que

$$P(Y > 6) = 1 - F(6) = e^{-3/2} = 0.2231.$$

Por lo tanto, la probabilidad de que la lavadora requiera una reparación mayor después de seis años es de 0.223. Desde luego, la probabilidad de que requiera reparación antes del sexto año es de 0.777. Así, se podría concluir que la lavadora no es realmente una ganga. La probabilidad de que se requiera una reparación mayor durante el primer año es

$$P(Y < 1) = 1 - e^{-1/4} = 1 - 0.779 = 0.221. \quad \text{J}$$

6.7 Distribución chi cuadrada

Otro caso especial muy importante de la distribución gamma se obtiene al permitir que $\alpha = \nu/2$ y $\beta = 2$, donde ν es un entero positivo. Este resultado se conoce como **distribución chi cuadrada**. La distribución tiene un solo parámetro, ν , denominado **grados de libertad**.

Distribución chi cuadrada La variable aleatoria continua X tiene una **distribución chi cuadrada**, con ν **grados de libertad**, si su función de densidad es dada por

$$f(x; \nu) = \begin{cases} \frac{1}{2^{\nu/2} \Gamma(\nu/2)} x^{\nu/2-1} e^{-x/2}, & x > 0, \\ 0, & \text{en otro caso,} \end{cases}$$

donde ν es un entero positivo.

La distribución chi cuadrada desempeña un papel fundamental en la inferencia estadística. Tiene una aplicación considerable tanto en la metodología como en la teoría. Aunque no estudiaremos con detalle sus aplicaciones en este capítulo, es importante tener en cuenta que los capítulos 8, 9 y 16 contienen aplicaciones importantes. La distribución chi cuadrada es un componente importante de la prueba estadística de hipótesis y de la estimación estadística.

Los temas en los que se trata con distribuciones de muestreo, análisis de varianza y estadística no paramétrica implican el uso extenso de la distribución chi cuadrada.

Teorema 6.5: La media y la varianza de la distribución chi cuadrada son

$$\mu = \nu \text{ y } \sigma^2 = 2\nu.$$

6.8 Distribución beta

Una extensión de la distribución uniforme es la distribución beta. Primero definiremos una **función beta**.

Definición 6.3: Una **función beta** es definida por

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}, \text{ para } \alpha, \beta > 0,$$

donde $\Gamma(\alpha)$ es la función gamma.

Distribución beta La variable aleatoria continua X tiene una **distribución beta** con los parámetros $\alpha > 0$ y $\beta > 0$, si su función de densidad es dada por

$$f(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 < x < 1, \\ 0, & \text{en otro caso.} \end{cases}$$

Observe que la distribución uniforme sobre $(0, 1)$ es una distribución beta con los parámetros $\alpha = 1$ y $\beta = 1$.

Teorema 6.6: La media y la varianza de una distribución beta en la que los parámetros α y β son

$$\mu = \frac{\alpha}{\alpha + \beta} \text{ y } \sigma^2 = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)},$$

respectivamente.

Para la distribución uniforme sobre $(0, 1)$, la media y la varianza son

$$\mu = \frac{1}{1+1} = \frac{1}{2} \text{ y } \sigma^2 = \frac{(1)(1)}{(1+1)^2(1+1+1)} = \frac{1}{12},$$

respectivamente.

6.9 Distribución logarítmica normal

La distribución logarítmica normal se utiliza en una amplia variedad de aplicaciones. La distribución se aplica en casos donde una transformación logarítmica natural tiene como resultado una distribución normal.

Distribución logarítmica normal La variable aleatoria continua X tiene una **distribución logarítmica normal** si la variable aleatoria $Y = \ln(X)$ tiene una distribución normal con media μ y desviación estándar σ . La función de densidad de X que resulta es

$$f(x; \mu, \sigma) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(\ln(x)-\mu)^2}, & x \geq 0, \\ 0, & x < 0. \end{cases}$$

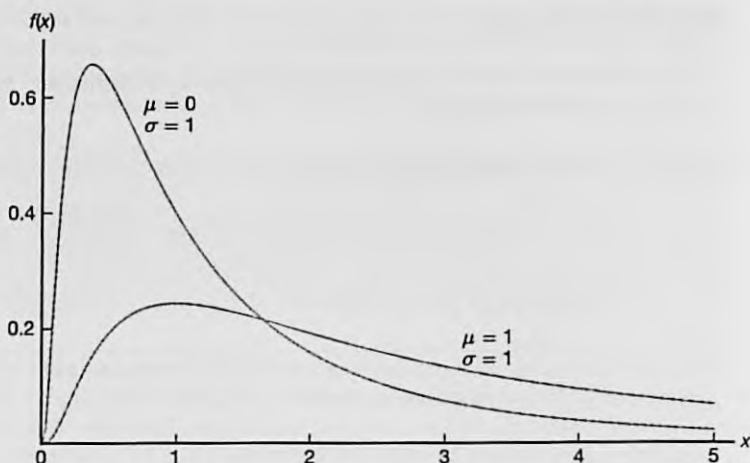


Figura 6.29: Distribuciones logarítmicas normales.

Las gráficas de las distribuciones logarítmicas normales se ilustran en la figura 6.29.

Teorema 6.7: La media y la varianza de la distribución logarítmica normal son

$$\mu = e^{\mu + \sigma^2/2} \text{ y } \sigma^2 = e^{2\mu + \sigma^2} (e^{\sigma^2} - 1).$$

La función de distribución acumulativa es muy simple debido a su relación con la distribución normal. El uso de la función de distribución se ilustra con el siguiente ejemplo.

Ejemplo 6.22: Se sabe que históricamente la concentración de contaminantes producidos por plantas químicas exhiben un comportamiento que se parece a una distribución logarítmica normal. Esto es importante cuando se consideran cuestiones relacionadas con el cumplimiento de las regulaciones gubernamentales. Suponga que la concentración de cierto contaminante, en partes por millón, tiene una distribución logarítmica normal con los parámetros $\mu = 3.2$ y $\sigma = 1$. ¿Cuál es la probabilidad de que la concentración exceda 8 partes por millón?

Solución: Sea la variable aleatoria X la concentración de contaminantes. Entonces

$$P(X > 8) = 1 - P(X \leq 8).$$

Como $\ln(X)$ tiene una distribución normal con media $\mu = 3.2$ y desviación estándar $\sigma = 1$,

$$P(X \leq 8) = \Phi \left[\frac{\ln(8) - 3.2}{1} \right] = \Phi(-1.12) = 0.1314.$$

Aquí, utilizamos el símbolo Φ para denotar la función de distribución acumulativa de la distribución normal estándar. Como resultado, la probabilidad de que la concentración del contaminante exceda 8 partes por millón es 0.1314. J

Ejemplo 6.23: La vida, en miles de millas, de un cierto tipo de control electrónico para locomotoras tiene una distribución aproximadamente logarítmica normal con $\mu = 5.149$ y $\sigma = 0.737$. Calcule el quinto percentil de la vida de un control electrónico como éste.

Solución: A partir de la tabla A.3 sabemos que $P(Z < -1.645) = 0.05$. Denote como X la vida del control electrónico. Puesto que $\ln(X)$ tiene una distribución normal con media $\mu = 5.149$ y $\sigma = 0.737$, el quinto percentil de X se calcula como

$$\ln(x) = 5.149 + (0.737)(-1.645) = 3.937.$$

Por lo tanto, $x = 51.265$. Esto significa que sólo 5% de los controles tendrán un tiempo de vida menor que 51,265 millas. ■

6.10 Distribución de Weibull (opcional)

La tecnología actual permite que los ingenieros diseñen muchos sistemas complicados cuya operación y seguridad dependen de la confiabilidad de los diversos componentes que conforman los sistemas. Por ejemplo, un fusible se puede quemar, una columna de acero se puede torcer o un dispositivo sensor de calor puede fallar. Componentes idénticos, sujetos a idénticas condiciones ambientales, fallarán en momentos diferentes e impredecibles. Ya examinamos el papel que desempeñan las distribuciones gamma y exponencial en estos tipos de problemas. Otra distribución que se ha utilizado ampliamente en años recientes para tratar con tales problemas es la **distribución de Weibull**, introducida por el físico sueco Waloddi Weibull en 1939.

Distribución de Weibull La variable aleatoria continua X tiene una **distribución de Weibull**, con parámetros α y β , si su función de densidad es dada por

$$f(x; \alpha, \beta) = \begin{cases} \alpha\beta x^{\beta-1} e^{-\alpha x^\beta}, & x > 0, \\ 0, & \text{en otro caso,} \end{cases}$$

donde $\alpha > 0$ y $\beta > 0$.

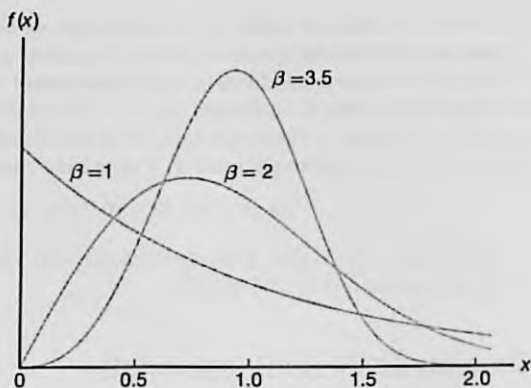
En la figura 6.30 se ilustran las gráficas de la distribución de Weibull para $\alpha = 1$ y diversos valores del parámetro β . Vemos que las curvas cambian de manera considerable para diferentes valores del parámetro β . Si permitimos que $\beta = 1$, la distribución de Weibull se reduce a la distribución exponencial. Para valores de $\beta > 1$ las curvas adoptan ligeramente la forma de campana y se asemejan a las curvas normales, pero muestran algo de asimetría.

La media y la varianza de la distribución de Weibull se establecen en el siguiente teorema. Se solicita al lector que haga la demostración en el ejercicio 6.52 de la página 206.

Teorema 6.8: La media y la varianza de la distribución de Weibull son

$$\mu = \alpha^{-1/\beta} \Gamma\left(1 + \frac{1}{\beta}\right) \quad \text{y} \quad \sigma^2 = \alpha^{-2/\beta} \left\{ \Gamma\left(1 + \frac{2}{\beta}\right) - \left[\Gamma\left(1 + \frac{1}{\beta}\right) \right]^2 \right\}.$$

Al igual que la distribución gamma y la exponencial, la distribución de Weibull se aplica a problemas de confiabilidad y de prueba de vida como los de **tiempo de operación**

Figura 6.30: Distribuciones de Weibull ($\alpha = 1$).

antes de la falla o la duración de la vida de un componente, que se miden desde algún tiempo específico hasta que falla. Representemos este tiempo de operación antes de la falla mediante la variable aleatoria continua T , con función de densidad de probabilidad $f(t)$, donde $f(t)$ es la distribución de Weibull. Ésta tiene la flexibilidad inherente de no requerir la propiedad de falta de memoria de la distribución exponencial. La función de distribución acumulativa (fda) para la distribución de Weibull se puede escribir en forma cerrada y realmente es muy útil para calcular probabilidades.

Fda para la distribución de Weibull La función de distribución acumulativa para la distribución de Weibull es dada por

$$F(x) = 1 - e^{-\alpha x^\beta}, \quad \text{para } x \geq 0,$$

para $\alpha > 0$ y $\beta > 0$.

Ejemplo 6.24: El tiempo de vida X , en horas, de un artículo en el taller mecánico tiene una distribución de Weibull con $\alpha = 0.01$ y $\beta = 2$. ¿Cuál es la probabilidad de que falle antes de 8 horas de uso?

Solución: $P(X < 8) = F(8) = 1 - e^{-(0.01)8^2} = 1 - 0.527 = 0.473$.

La tasa de fallas para la distribución de Weibull

Cuando se aplica la distribución de Weibull, con frecuencia es útil determinar la **tasa de fallas** (algunas veces denominada tasa de riesgo) para tener conocimiento del desgaste o deterioro del componente. Comencemos por definir la confiabilidad de un componente o producto como la *probabilidad de que funcione adecuadamente por al menos un tiempo específico en condiciones experimentales específicas*. Por lo tanto, si $R(t)$ se define como la confiabilidad del componente dado en el tiempo t , escribimos

$$R(t) = P(T > t) = \int_t^\infty f(t) dt = 1 - F(t),$$

donde $F(t)$ es la función de distribución acumulativa de T . La probabilidad condicional de que un componente fallará en el intervalo de $T = t$ a $T = t + \Delta t$, dado que sobrevive hasta el tiempo t , es

$$\frac{F(t + \Delta t) - F(t)}{R(t)}$$

Al dividir esta proporción entre Δt y tomar el límite como $\Delta t \rightarrow 0$, obtenemos la tasa de fallas, denotada por $Z(t)$. De aquí,

$$Z(t) = \lim_{\Delta t \rightarrow 0} \frac{F(t + \Delta t) - F(t)}{\Delta t} \frac{1}{R(t)} = \frac{F'(t)}{R(t)} = \frac{f(t)}{R(t)} = \frac{f(t)}{1 - F(t)},$$

que expresa la tasa de fallas en términos de la distribución del tiempo de operación antes de la falla.

Como $Z(t) = f(t)/[1 - F(t)]$, entonces la tasa de falla es dada como sigue:

Tasa de fallas para la distribución de Weibull La **tasa de fallas** en el tiempo t para la distribución de Weibull es dada por

$$Z(t) = \alpha\beta t^{\beta-1}, \quad t > 0.$$

Interpretación de la tasa de fallas

La cantidad $Z(t)$ es bien llamada tasa de fallas porque realmente cuantifica la tasa de cambio con el tiempo de la probabilidad condicional de que el componente dure una Δt adicional *dado que ha durado el tiempo t* . La tasa de disminución (o crecimiento) con el tiempo también es importante. Los siguientes puntos son fundamentales.

- Si $\beta = 1$, la tasa de fallas $= \alpha$, es decir, una constante. Esto, como se indicó anteriormente, es el caso especial de la distribución exponencial en que predomina la falta de memoria.
- Si $\beta > 1$, $Z(t)$ es una función creciente del tiempo t que indica que el componente se desgasta con el tiempo.
- Si $\beta < 1$, $Z(t)$ es una función decreciente del tiempo t y, por lo tanto, el componente se fortalece o endurece con el paso del tiempo.

Por ejemplo, el artículo en el taller mecánico del ejemplo 6.24 tiene $\beta = 2$ y, por consiguiente, se desgasta con el tiempo. De hecho, la función de la tasa de fallas es dada por $Z(t) = .02t$. Por otro lado, suponga un parámetro donde $\beta = 3/4$ y $\alpha = 2$. En ese caso, $Z(t) = 1.5/t^{1/4}$ y, por lo tanto, el componente se hace más fuerte con el tiempo.

Ejercicios

6.39 Utilice la función gamma con $y = \sqrt{2x}$ para demostrar que $\Gamma(1/2) = \sqrt{\pi}$.

6.40 En cierta ciudad, el consumo diario de agua (en millones de litros) sigue aproximadamente una distribución gamma con $\alpha = 2$ y $\beta = 3$. Si la capacidad diaria de dicha ciudad es de 9 millones de litros de agua, ¿cuál es la probabilidad de que en cualquier día dado el suministro de agua sea inadecuado?

6.41 Si una variable aleatoria X tiene una distribución gamma con $\alpha = 2$ y $\beta = 1$, calcule $P(1.8 < X < 2.4)$.

6.42 Suponga que el tiempo, en horas, necesario para reparar una bomba de calor es una variable aleatoria X que tiene una distribución gamma con los parámetros $\alpha = 2$ y $\beta = 1/2$. ¿Cuál es la probabilidad de que la siguiente llamada de servicio requiera

- a lo sumo una hora para reparar la bomba de calor?
- al menos dos horas para reparar la bomba de calor?

6.43 a) Calcule la media y la varianza del consumo diario de agua del ejercicio 6.40.

- De acuerdo con el teorema de Chebyshev, ¿hay por lo menos $3/4$ de probabilidad de que el consumo de agua en cualquier día determinado caiga dentro de cuál intervalo?

6.44 En cierta ciudad el consumo diario de energía eléctrica, en millones de kilowatts-hora, es una variable aleatoria X que tiene una distribución gamma con media $\mu = 6$ y varianza $\sigma^2 = 12$.

- Calcule los valores de α y β .
- Calcule la probabilidad de que en cualquier día dado el consumo diario de energía exceda los 12 millones de kilowatts-hora.

6.45 El tiempo necesario para que un individuo sea atendido en una cafetería es una variable aleatoria que tiene una distribución exponencial con una media de 4 minutos. ¿Cuál es la probabilidad de que una persona sea atendida en menos de 3 minutos en al menos 4 de los siguientes 6 días?

6.46 La vida, en años, de cierto interruptor eléctrico tiene una distribución exponencial con una vida promedio de $\beta = 2$. Si 100 de estos interruptores se instalan en diferentes sistemas, ¿cuál es la probabilidad de que, a lo sumo, fallen 30 durante el primer año?

6.47 Suponga que la vida de servicio de la batería de un auxiliar auditivo, en años, es una variable aleatoria que tiene una distribución de Weibull con $\alpha = 1/2$ y $\beta = 2$.

- ¿Cuánto tiempo se puede esperar que dure tal batería?
- ¿Cuál es la probabilidad de que tal batería esté funcionando después de 2 años?

6.48 Derive la media y la varianza de la distribución beta.

6.49 Suponga que la variable aleatoria X tiene una distribución beta con $\alpha = 1$ y $\beta = 3$.

- Determine la media y la mediana de X .
- Determine la varianza de X .
- Calcule la probabilidad de que $X > 1/3$.

6.50 Si la proporción de una marca de televisores que requiere servicio durante el primer año de operación es una variable aleatoria que tiene una distribución beta con $\alpha = 3$ y $\beta = 2$, ¿cuál es la probabilidad de que al menos 80% de los nuevos modelos de esta marca que se vendieron este año requieran servicio durante su primer año de operación?

6.51 Las vidas de ciertos sellos para automóvil tienen la distribución de Weibull con tasa de fallas $Z(t) = 1/\sqrt{t}$. Calcule la probabilidad de que tal sello aún esté intacto después de 4 años.

6.52 Derive la media y la varianza de la distribución de Weibull.

6.53 En una investigación biomédica se determinó que el tiempo de supervivencia, en semanas, de un animal cuando se le somete a cierta exposición de radiación gamma tiene una distribución gamma con $\alpha = 5$ y $\beta = 10$.

- ¿Cuál es el tiempo medio de supervivencia de un animal seleccionado al azar del tipo que se utilizó en el experimento?
- ¿Cuál es la desviación estándar del tiempo de supervivencia?
- ¿Cuál es la probabilidad de que un animal sobreviva más de 30 semanas?

6.54 Se sabe que la vida, en semanas, de cierto tipo de transistor tiene una distribución gamma con una media de 10 semanas y una desviación estándar de $\sqrt{50}$ semanas.

- ¿Cuál es la probabilidad de que un transistor de este tipo dure a lo sumo 50 semanas?
- ¿Cuál es la probabilidad de que un transistor de este tipo no sobreviva las primeras 10 semanas?

6.55 El tiempo de respuesta de una computadora es una aplicación importante de las distribuciones gamma y exponencial. Suponga que un estudio de cierto sistema de cómputo revela que el tiempo de respuesta, en segundos, tiene una distribución exponencial con una media de 3 segundos.

- a) ¿Cuál es la probabilidad de que el tiempo de respuesta exceda 5 segundos?
 b) ¿Cuál es la probabilidad de que el tiempo de respuesta exceda 10 segundos?

6.56 Los datos de frecuencia a menudo tienen una distribución logarítmica normal. Se estudia el uso promedio de potencia (dB por hora) para una empresa específica y se sabe que tiene una distribución logarítmica normal con parámetros $\mu = 4$ y $\sigma = 2$. ¿Cuál es la probabilidad de que la empresa utilice más de 270 dB durante cualquier hora particular?

6.57 Para el ejercicio 6.56, ¿cuál es el uso de la potencia media (dB promedio por hora)? ¿Cuál es la varianza?

6.58 El número de automóviles que llegan a cierta intersección por minuto tiene una distribución de Poisson con una media de 5. Existe interés por el tiempo que transcurre antes de que 10 automóviles aparezcan en la intersección.

- a) ¿Cuál es la probabilidad de que más de 10 automóviles aparezcan en la intersección durante cualquier minuto determinado?
 b) ¿Cuál es la probabilidad de que transcurran más de 2 minutos antes de que lleguen 10 autos?

6.59 Considere la información del ejercicio 6.58.

- a) ¿Cuál es la probabilidad de que transcurra más de 1 minuto entre llegadas?
 b) ¿Cuál es el número medio de minutos que transcurre entre las llegadas?

6.60 Demuestre que la función de la tasa de fallas es dada por

$$Z(t) = \alpha \beta t^{\beta-1}, \quad t > 0,$$

si y sólo si la distribución del tiempo que transcurre antes de la falla es la distribución de Weibull

$$f(t) = \alpha \beta t^{\beta-1} e^{-\alpha t^\beta}, \quad t > 0.$$

Ejercicios de repaso

6.61 Según un estudio publicado por un grupo de sociólogos de la Universidad de Massachusetts, aproximadamente 49% de los consumidores de Valium en el estado de Massachusetts son empleados de oficina. ¿Cuál es la probabilidad de que entre 482 y 510 de los siguientes 1000 consumidores de Valium seleccionados al azar de dicho estado sean empleados de oficina?

6.62 La distribución exponencial se aplica con frecuencia a los tiempos de espera entre éxitos en un proceso de Poisson. Si el número de llamadas que se reciben por hora en un servicio de respuesta telefónica es una variable aleatoria de Poisson con el parámetro $\lambda = 6$, sabemos que el tiempo, en horas, entre llamadas sucesivas tiene una distribución exponencial con el parámetro $\beta = 1/6$. ¿Cuál es la probabilidad de esperar más de 15 minutos entre cualesquiera 2 llamadas sucesivas?

6.63 Cuando α es un entero positivo n , la distribución gamma también se conoce como **distribución de Erlang**. Al establecer que $\alpha = n$ en la distribución gamma de la página 195, la distribución de Erlang es

$$f(x) = \begin{cases} \frac{x^{n-1} e^{-x/\beta}}{\beta^n (n-1)!}, & x > 0, \\ 0, & \text{en otro caso.} \end{cases}$$

Se puede demostrar que si los tiempos entre eventos sucesivos son independientes, y cada uno tiene una distribución exponencial con el parámetro β , entonces el tiempo de espera total X transcurrido hasta que ocurran n eventos tiene la distribución de Erlang. Con referen-

cia al ejercicio de repaso 6.62, ¿cuál es la probabilidad de que las siguientes 3 llamadas se reciban dentro de los siguientes 30 minutos?

6.64 Un fabricante de cierto tipo de máquina grande desea comprar remaches de uno de dos fabricantes. Es importante que la resistencia a la rotura de cada remache exceda 10,000 psi. Dos fabricantes (A y B) ofrecen este tipo de remache y ambos tienen remaches cuya resistencia a la rotura está distribuida de forma normal. Las resistencias promedio a la rotura para los fabricantes A y B son 14,000 psi y 13,000 psi, respectivamente. Las desviaciones estándar son 2000 psi y 1000 psi, respectivamente. ¿Cuál fabricante producirá, en promedio, el menor número de remaches defectuosos?

6.65 De acuerdo con un censo reciente, casi 65% de los hogares en Estados Unidos se componen de una o dos personas. Si se supone que este porcentaje sigue siendo válido en la actualidad, ¿cuál es la probabilidad de que entre 590 y 625 de los siguientes 1000 hogares seleccionados al azar en Estados Unidos consten de una o dos personas?

6.66 Cierta tipo de dispositivo tiene una tasa de fallas anunciada de 0.01 por hora. La tasa de fallas es constante y se aplica la distribución exponencial.

- a) ¿Cuál es el tiempo promedio que transcurre antes de la falla?
 b) ¿Cuál es la probabilidad de que pasen 200 horas antes de que se observe una falla?

6.67 En una planta de procesamiento químico es importante que el rendimiento de cierto tipo de producto de un lote se mantenga por arriba de 80%. Si permanece por debajo de 80% durante un tiempo prolongado, la empresa pierde dinero. Los lotes producidos ocasionalmente con defectos son de poco interés, pero si varios lotes por día resultan defectuosos, la planta se detiene y se realizan ajustes. Se sabe que el rendimiento se distribuye normalmente con una desviación estándar de 4%.

- ¿Cuál es la probabilidad de una "falsa alarma" (rendimiento por debajo de 80%) cuando el rendimiento promedio es en realidad de 85%?
- ¿Cuál es la probabilidad de que un lote tenga un rendimiento que exceda el 80% cuando en realidad el rendimiento promedio es de 79%?

6.68 Para un componente eléctrico que tiene una tasa de fallas de una vez cada 5 horas es importante considerar el tiempo que transcurre para que fallen 2 componentes.

- Suponiendo que se aplica la distribución gamma, ¿cuál es el tiempo promedio que transcurre para que fallen 2 componentes?
- ¿Cuál es la probabilidad de que transcurran 12 horas antes de que fallen 2 componentes?

6.69 Se establece que la elongación de una barra de acero bajo una carga particular se distribuye normalmente con una media de 0.05 pulgadas y $\sigma = 0.01$ pulgadas. Calcule la probabilidad de que el alargamiento esté

- por arriba de 0.1 pulgadas;
- por abajo de 0.04 pulgadas;
- entre 0.025 y 0.065 pulgadas.

6.70 Se sabe que un satélite controlado tiene un error (distancia del objetivo) que se distribuye normalmente con una media 0 y una desviación estándar de 4 pies. El fabricante del satélite define un éxito como un disparo en el cual el satélite llega a 10 pies del objetivo. Calcule la probabilidad de que el satélite falle.

6.71 Un técnico planea probar cierto tipo de resina desarrollada en el laboratorio para determinar la naturaleza del tiempo que transcurre antes de que se logre el pegado. Se sabe que el tiempo promedio para el pegado es de 3 horas y que la desviación estándar es de 0.5 horas. Un producto se considerará indeseable si el tiempo de pegado es menor de una hora o mayor de 4 horas. Comente sobre la utilidad de la resina. ¿Con qué frecuencia su desempeño se considera indeseable? Suponga que el tiempo para la unión se distribuye normalmente.

6.72 Considere la información del ejercicio de repaso 6.66. ¿Cuál es la probabilidad de que transcurran menos de 200 horas antes de que ocurran 2 fallas?

6.73 Para el ejercicio de repaso 6.72, ¿cuál es la media y la varianza del tiempo que transcurre antes de que ocurran 2 fallas?

6.74 Se sabe que la tasa promedio de uso de agua (en miles de galones por hora) en cierta comunidad implica la distribución logarítmica normal con los parámetros $\mu = 5$ y $\sigma = 2$. Para propósitos de planeación es importante tener información sobre los períodos de alto consumo. ¿Cuál es la probabilidad de que, para cualquier hora determinada, se usen 50,000 galones de agua?

6.75 Para el ejercicio de repaso 6.74, ¿cuál es la media del uso de agua por hora promedio en miles de galones?

6.76 En el ejercicio 6.54 de la página 206 se supone que la vida de un transistor tiene una distribución gamma con una media de 10 semanas y una desviación estándar de $\sqrt{50}$ semanas. Suponga que la distribución gamma es incorrecta y que se trata de una distribución normal.

- ¿Cuál es la probabilidad de que el transistor dure a lo sumo 50 semanas?
- ¿Cuál es la probabilidad de que el transistor no sobreviva las primeras 10 semanas?
- Comente acerca de la diferencia entre los resultados que obtuvo aquí y los que se obtuvieron en el ejercicio 6.54 de la página 206.

6.77 La distribución beta tiene muchas aplicaciones en problemas de confiabilidad, donde la variable aleatoria básica es una proporción, como sucede en el contexto práctico que se ilustra en el ejercicio 6.50 de la página 206. En este apartado considere el ejercicio de repaso 3.73 de la página 108. Las impurezas en el lote del producto de un proceso químico reflejan un problema grave. Se sabe que la proporción de impurezas Y en un lote tiene la siguiente función de densidad

$$f(y) = \begin{cases} 10(1-y)^9, & 0 \leq y \leq 1, \\ 0, & \text{en otro caso.} \end{cases}$$

- Verifique que la anterior sea una función de densidad válida.
- ¿Cuál es la probabilidad de que un lote se considere no aceptable (es decir, $Y > 0.6$)?
- ¿Cuáles son los parámetros α y β de la distribución beta que se ilustra aquí?
- La media de la distribución beta es $\frac{\alpha}{\alpha+\beta}$. ¿Cuál es la proporción media de impurezas en el lote?
- La varianza de una variable aleatoria beta distribuida es

$$\sigma^2 = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}.$$

¿Cuál es la varianza de Y en este problema?

6.78 Considere ahora el ejercicio de repaso 3.74 de la página 108. La función de densidad del tiempo Z entre las llamadas, en minutos, a una empresa de suministro eléctrico es dada por

$$f(z) = \begin{cases} \frac{1}{10} e^{-z/10}, & 0 < z < \infty, \\ 0, & \text{en otro caso.} \end{cases}$$

- ¿Cuál es el tiempo medio entre llamadas?
- ¿Cuál es la varianza en el tiempo entre llamadas?
- ¿Cuál es la probabilidad de que el tiempo entre llamadas supere la media?

6.79 Considere el ejercicio de repaso 6.78. Dada la suposición de la distribución exponencial, ¿cuál es el número medio de llamadas por hora? ¿Cuál es la varianza en el número de llamadas por hora?

6.80 En un proyecto experimental sobre el factor humano se determinó que el tiempo de reacción de un piloto ante un estímulo visual es distribuido normalmente con una media de $1/2$ segundo y una desviación estándar de $2/5$ de segundo.

- ¿Cuál es la probabilidad de que una reacción del piloto tome más de 0.3 segundos?
- ¿Qué tiempo de reacción se excede el 95% de las veces?

6.81 El tiempo que transcurre entre las fallas de una pieza esencial de equipo es importante en la decisión del uso de equipo auxiliar. Un ingeniero cree que el mejor modelo para el tiempo entre las fallas de un generador es la distribución exponencial con una media de 15 días.

- Si el generador acaba de fallar, ¿cuál es la probabilidad de que falle en los siguientes 21 días?
- ¿Cuál es la probabilidad de que el generador funcione durante 30 días sin fallar?

6.82 El periodo de vida de una broca en una operación mecánica, en horas, tiene una distribución de Weibull con $\alpha = 2$ y $\beta = 50$. Calcule la probabilidad de que la broca falle antes de 10 horas de uso.

6.83 Calcule la fda para la distribución de Weibull. [Sugerencia: En la definición de una fda haga la transformación $z = y^\beta$].

6.84 Explique por qué la naturaleza del escenario en el ejercicio de repaso 6.82 probablemente no se preste a la distribución exponencial.

6.85 A partir de la relación entre la variable aleatoria chi cuadrada y la variable aleatoria gamma, demuestre que la media de la variable aleatoria chi cuadrada es v y que la varianza es $2v$.

6.86 El tiempo que le toma a un usuario de computadora leer su correo electrónico, en segundos, se distribuye como una variable aleatoria logarítmica normal con $\mu = 1.8$ y $\sigma^2 = 4.0$.

- ¿Cuál es la probabilidad de que el usuario lea el correo durante más de 20 segundos? ¿Y por más de un minuto?
- ¿Cuál es la probabilidad de que el usuario lea el correo durante un tiempo que sea igual a la media de la distribución logarítmica normal subyacente?

6.87 Proyecto de grupo: Pida a grupos de estudiantes que observen durante 2 semanas el número de personas que entra a una cafetería o restaurante de comida rápida específico en el transcurso de una hora, empezando a la misma hora cada día. La hora deberá ser la de mayor tránsito en la cafetería o restaurante. Los datos reunidos corresponderán al número de clientes que entran al lugar durante cada lapso de media hora. De esta manera, cada día se recolectarán 2 datos. Supongamos que la variable aleatoria X , el número de personas que entra cada media hora, tiene una distribución de Poisson. Los estudiantes deberán calcular la media y la varianza muestrales de X utilizando los 28 datos obtenidos.

- ¿Qué evidencia hay de que la distribución de Poisson es o no correcta?
- Dado que X es una variable de Poisson, ¿cuál es la distribución de T , el tiempo entre la llegada de las personas al lugar durante un lapso de media hora? Proporcione un estimado numérico del parámetro de esa distribución.
- Proporcione un estimado de la probabilidad de que el lapso de tiempo entre las 2 llegadas sea menor de 15 minutos.
- ¿Cuál es la probabilidad estimada de que el lapso entre las 2 llegadas sea mayor de 10 minutos?
- ¿Cuál es la probabilidad estimada de que 20 minutos después de iniciar la recolección de datos ningún cliente haya llegado?

6.11 Posibles riesgos y errores conceptuales; relación con el material de otros capítulos

Muchos de los riesgos en el uso del material de este capítulo son muy similares a los del capítulo 5. Uno de los peores abusos de la estadística consiste en suponer que se trata de

una distribución normal haciendo algún tipo de inferencia estadística, cuando en realidad no es normal. En los capítulos 10 al 15 el lector estudiará las pruebas de hipótesis, en las que se asume normalidad. Además, se le recordará al lector que hay **pruebas de la bondad de ajuste**, además de las rutinas gráficas que se examinan en los capítulos 8 y 10, que permiten verificar los datos para determinar si es razonable la suposición de normalidad.

Debemos hacer advertencias similares con respecto a las suposiciones que a menudo se hacen sobre otras distribuciones, además de la curva normal. En este libro se han presentado ejemplos en los que es necesario calcular las probabilidades de falla de ciertos productos o la probabilidad de recibir una queja durante cierto periodo. Se suelen hacer suposiciones con respecto a cierto tipo de distribución, así como a los valores de los parámetros de la distribución. Observe que los problemas de ejemplo incluyen los valores de los parámetros (por ejemplo, el valor de β para la distribución exponencial). No obstante, en los problemas de la vida real los valores de los parámetros deben ser estimaciones de experiencias o datos reales. Observe el énfasis que se pone en la estimación en los proyectos que aparecen en los capítulos 1, 5 y 6, así como la referencia que se hace en el capítulo 5 a las estimación de parámetros, tema que se analizará ampliamente a partir del capítulo 9.

Capítulo 7

Funciones de variables aleatorias (opcional)

7.1 Introducción

Este capítulo contiene un amplio espectro de material. Los capítulos 5 y 6 tratan tipos específicos de distribuciones, tanto discretas como continuas. Éstas son distribuciones que suelen aplicarse en muchos campos, por ejemplo en el de la confiabilidad, el de control de calidad y el de muestreo de aceptación. En este capítulo comenzamos a estudiar un tema más general: el de la distribución de funciones de variables aleatorias. Se presentan las técnicas generales y se ilustran con ejemplos. Las presentaciones van seguidas por un concepto relacionado, el de *funciones generadoras de momentos*, que pueden ser útiles para el aprendizaje de distribuciones de funciones lineales de variables aleatorias.

En los métodos estadísticos estándar, el resultado de la prueba de hipótesis estadísticas, la estimación, o incluso las gráficas estadísticas, no involucra a una sola variable aleatoria sino a *funciones de una o más variables aleatorias*. Como resultado, la inferencia estadística requiere la distribución de tales funciones. Por ejemplo, es común que se utilicen **promedios de variables aleatorias**. Además, las sumatorias y las combinaciones lineales más generales son importantes. Con frecuencia nos interesa la distribución de las sumas de cuadrados de variables aleatorias, en particular la manera en que se utilizan las técnicas del análisis de varianza, las cuales se estudiarán en los capítulos 11 a 14.

7.2 Transformaciones de variables

Con frecuencia, en la estadística se enfrenta la necesidad de derivar la distribución de probabilidad de una función de una o más variables aleatorias. Por ejemplo, suponga que X es una variable aleatoria discreta con distribución de probabilidad $f(x)$, suponga también que $Y = u(X)$ define una transformación uno a uno entre los valores de X y Y . Queremos encontrar la distribución de probabilidad de Y . Es importante notar que la transformación uno a uno implica que cada valor x está relacionado con un, y sólo un, valor $y = u(x)$, y que cada valor y está relacionado con un, y sólo un, valor $x = w(y)$, donde $w(y)$ se obtiene al resolver $y = u(x)$ para x en términos de y .

A partir de lo expuesto respecto a las distribuciones de probabilidad discreta en el capítulo 3, nos quedó claro que la variable aleatoria Y toma el valor y cuando X toma el valor $w(y)$. En consecuencia, la distribución de probabilidad de Y es dada por

$$g(y) = P(Y = y) = P[X = w(y)] = f[w(y)].$$

Teorema 7.1: Suponga que X es una variable aleatoria **discreta** con distribución de probabilidad $f(x)$. Definamos con $Y = u(X)$ una transformación uno a uno entre los valores de X y Y , de manera que la ecuación $y = u(x)$ se resuelva exclusivamente para x en términos de y , digamos, $x = w(y)$. Entonces, la distribución de probabilidad de Y es

$$g(y) = f[w(y)].$$

Ejemplo 7.1: Sea X una variable aleatoria geométrica con la siguiente distribución de probabilidad

$$f(x) = \frac{3}{4} \left(\frac{1}{4}\right)^{x-1}, \quad x = 1, 2, 3, \dots$$

Calcule la distribución de probabilidad de la variable aleatoria $Y = X^2$.

Solución: Como todos los valores de X son positivos, la transformación define una correspondencia uno a uno entre los valores x y y , $y = x^2$ y $x = \sqrt{y}$. Por lo tanto,

$$g(y) = \begin{cases} f(\sqrt{y}) = \frac{3}{4} \left(\frac{1}{4}\right)^{\sqrt{y}-1}, & y = 1, 4, 9, \dots \\ 0, & \text{en cualquier caso.} \end{cases}$$

De manera similar, para una transformación de dos dimensiones, tenemos el resultado en el teorema 7.2.

Teorema 7.2: Suponga que X_1 y X_2 son variables aleatorias **discretas**, con distribución de probabilidad conjunta $f(x_1, x_2)$. Definamos con $Y_1 = u_1(X_1, X_2)$ y $Y_2 = u_2(X_1, X_2)$ una transformación uno a uno entre los puntos (x_1, x_2) y (y_1, y_2) , de manera que las ecuaciones

$$y_1 = u_1(x_1, x_2) \quad \text{y} \quad y_2 = u_2(x_1, x_2)$$

se pueden resolver exclusivamente para x_1 y x_2 en términos de y_1 y y_2 , digamos $x_1 = w_1(y_1, y_2)$ y $x_2 = w_2(y_1, y_2)$. Entonces, la distribución de probabilidad conjunta de Y_1 y Y_2 es

$$g(y_1, y_2) = f[w_1(y_1, y_2), w_2(y_1, y_2)].$$

El teorema 7.2 es muy útil para encontrar la distribución de alguna variable aleatoria $Y_1 = u_1(X_1, X_2)$, donde X_1 y X_2 son variables aleatorias discretas con distribución de probabilidad conjunta $f(x_1, x_2)$. Definimos simplemente una segunda función, digamos $Y_2 = u_2(X_1, X_2)$, manteniendo una correspondencia uno a uno entre los puntos (x_1, x_2) y (y_1, y_2) , y obtenemos la distribución de probabilidad conjunta $g(y_1, y_2)$. La distribución de Y_1 es precisamente la distribución marginal de $g(y_1, y_2)$ que se encuentra sumando los valores y_2 . Si denotamos la distribución de Y_1 con $h(y_1)$, podemos escribir

$$h(y_1) = \sum_{y_2} g(y_1, y_2).$$

Ejemplo 7.2: Sean X_1 y X_2 dos variables aleatorias independientes que tienen distribuciones de Poisson con los parámetros μ_1 y μ_2 , respectivamente. Calcule la distribución de la variable aleatoria $Y_1 = X_1 + X_2$.

Solución: Como X_1 y X_2 son independientes, podemos escribir

$$f(x_1, x_2) = f(x_1)f(x_2) = \frac{e^{-\mu_1} \mu_1^{x_1}}{x_1!} \frac{e^{-\mu_2} \mu_2^{x_2}}{x_2!} = \frac{e^{-(\mu_1 + \mu_2)} \mu_1^{x_1} \mu_2^{x_2}}{x_1! x_2!},$$

donde $x_1 = 0, 1, 2, \dots$ y $x_2 = 0, 1, 2, \dots$. Definamos ahora una segunda variable aleatoria, digamos $Y_2 = X_2$. Las funciones inversas son dadas por $x_1 = y_1 - y_2$ y $x_2 = y_2$. Si usamos el teorema 7.2, encontramos que la distribución de probabilidad conjunta de Y_1 y Y_2 es

$$g(y_1, y_2) = \frac{e^{-(\mu_1 + \mu_2)} \mu_1^{y_1 - y_2} \mu_2^{y_2}}{(y_1 - y_2)! y_2!},$$

donde $y_1 = 0, 1, 2, \dots$ y $y_2 = 0, 1, 2, \dots, y_1$. Advierta que, como $x_1 > 0$, la transformación $x_1 = y_1 - y_2$ implica que y_2 , por lo tanto, x_2 siempre deben ser menores o iguales que y_1 . En consecuencia, la distribución de probabilidad marginal de Y_1 es

$$\begin{aligned} h(y_1) &= \sum_{y_2=0}^{y_1} g(y_1, y_2) = e^{-(\mu_1 + \mu_2)} \sum_{y_2=0}^{y_1} \frac{\mu_1^{y_1 - y_2} \mu_2^{y_2}}{(y_1 - y_2)! y_2!} \\ &= \frac{e^{-(\mu_1 + \mu_2)}}{y_1!} \sum_{y_2=0}^{y_1} \frac{y_1!}{y_2! (y_1 - y_2)!} \mu_1^{y_1 - y_2} \mu_2^{y_2} \\ &= \frac{e^{-(\mu_1 + \mu_2)}}{y_1!} \sum_{y_2=0}^{y_1} \binom{y_1}{y_2} \mu_1^{y_1 - y_2} \mu_2^{y_2}. \end{aligned}$$

Al reconocer esta suma como la expansión binomial de $(\mu_1 + \mu_2)^{y_1}$, obtenemos

$$h(y_1) = \frac{e^{-(\mu_1 + \mu_2)} (\mu_1 + \mu_2)^{y_1}}{y_1!}, \quad y_1 = 0, 1, 2, \dots,$$

a partir de lo cual concluimos que la suma de las dos variables aleatorias independientes que tienen distribuciones de Poisson, con los parámetros μ_1 y μ_2 , tiene una distribución de Poisson con el parámetro $\mu_1 + \mu_2$. \blacksquare

Para calcular la distribución de probabilidad de la variable aleatoria $Y = u(X)$, cuando X es una variable aleatoria continua y la transformación es uno a uno, necesitaremos el teorema 7.3. La demostración de este teorema se deja al lector.

Teorema 7.3: Suponga que X es una variable aleatoria **continua** con distribución de probabilidad $f(x)$. Definamos con $Y = u(X)$ una correspondencia uno a uno entre los valores de X y Y , de manera que la ecuación $y = u(x)$ se resuelva exclusivamente para x en términos de y , digamos $x = w(y)$. Entonces, la distribución de probabilidad de Y es

$$g(y) = f[w(y)]|J|,$$

donde $J = w'(y)$ y se llama **jacobiano** de la transformación.

Ejemplo 7.3: Sea X una variable aleatoria continua con la siguiente distribución de probabilidad

$$f(x) = \begin{cases} \frac{x}{12}, & 1 < x < 5, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la distribución de probabilidad de la variable aleatoria $Y = 2X - 3$.

Solución: La solución inversa de $y = 2x - 3$ produce $x = (y + 3)/2$, de la que obtenemos $J = w'(y) = dx/dy = 1/2$. Por lo tanto, usando el teorema 7.3 encontramos que la función de densidad de Y es

$$g(y) = \begin{cases} \frac{(y+3)/2}{12} \left(\frac{1}{2}\right) = \frac{y+3}{48}, & -1 < y < 7, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Para calcular la distribución de probabilidad conjunta de las variables aleatorias $Y_1 = u_1(X_1, X_2)$ y $Y_2 = u_2(X_1, X_2)$, cuando X_1 y X_2 son continuas y la transformación es uno a uno, necesitamos un teorema adicional análogo al teorema 7.2, el cual establecemos sin demostración.

Teorema 7.4: Suponga que X_1 y X_2 son variables aleatorias **continuas** con distribución de probabilidad conjunta $f(x_1, x_2)$. Definamos con $Y_1 = u_1(X_1, X_2)$ y $Y_2 = u_2(X_1, X_2)$ una transformación uno a uno entre los puntos (x_1, x_2) y (y_1, y_2) , de manera que las ecuaciones $y_1 = u_1(x_1, x_2)$ y $y_2 = u_2(x_1, x_2)$ se resuelven exclusivamente para x_1 y x_2 en términos de y_1 y y_2 , digamos $x_1 = w_1(y_1, y_2)$ y $x_2 = w_2(y_1, y_2)$. Entonces, la distribución de probabilidad conjunta de Y_1 y Y_2 es

$$g(y_1, y_2) = f[w_1(y_1, y_2), w_2(y_1, y_2)]|J|,$$

donde el jacobiano es el determinante 2×2

$$J = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{vmatrix}$$

y $\frac{\partial x_1}{\partial y_1}$ es simplemente la derivada de $x_1 = w_1(y_1, y_2)$ respecto a y_1 , con y_2 constante, que en cálculo se denomina derivada parcial de x_1 respecto a y_1 . Las otras derivadas parciales se definen de manera similar.

Ejemplo 7.4: Sean X_1 y X_2 dos variables aleatorias continuas con la siguiente distribución de probabilidad conjunta

$$f(x_1, x_2) = \begin{cases} 4x_1x_2, & 0 < x_1 < 1, 0 < x_2 < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la distribución de probabilidad conjunta de $Y_1 = X_1^2$ y $Y_2 = X_1X_2$.

Solución: Las soluciones inversas de $y_1 = x_1^2$ y $y_2 = x_1x_2$ son $x_1 = \sqrt{y_1}$ y $x_2 = y_2/\sqrt{y_1}$, de las que obtenemos

$$J = \begin{vmatrix} 1/(2\sqrt{y_1}) & 0 \\ -y_2/2y_1^{3/2} & 1/\sqrt{y_1} \end{vmatrix} = \frac{1}{2y_1}.$$

Para determinar el conjunto B de puntos en el plano y_1, y_2 en el que se traza el conjunto A de puntos en el plano x_1, x_2 escribimos

$$x_1 = \sqrt{y_1} \quad y \quad x_2 = y_2 / \sqrt{y_1}.$$

Luego, al establecer $x_1 = 0, x_2 = 0, x_1 = 1$ y $x_2 = 1$, las fronteras del conjunto A se transforman en $y_1 = 0, y_2 = 0, y_1 = 1$ y $y_2 = \sqrt{y_1}$ o $y_2^2 = y_1$. Las dos regiones se ilustran en la figura 7.1. Al trazar el conjunto $A = \{(x_1, x_2) \mid 0 < x_1 < 1, 0 < x_2 < 1\}$ en el conjunto $B = \{(y_1, y_2) \mid y_2^2 < y_1 < 1, 0 < y_2 < 1\}$, se vuelve evidente que la transformación es uno a uno. Del teorema 7.4, la distribución de probabilidad conjunta de Y_1 y Y_2 es

$$g(y_1, y_2) = 4(\sqrt{y_1}) \frac{y_2}{\sqrt{y_1}} \frac{1}{2y_1} = \begin{cases} \frac{2y_2}{y_1}, & y_2^2 < y_1 < 1, 0 < y_2 < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

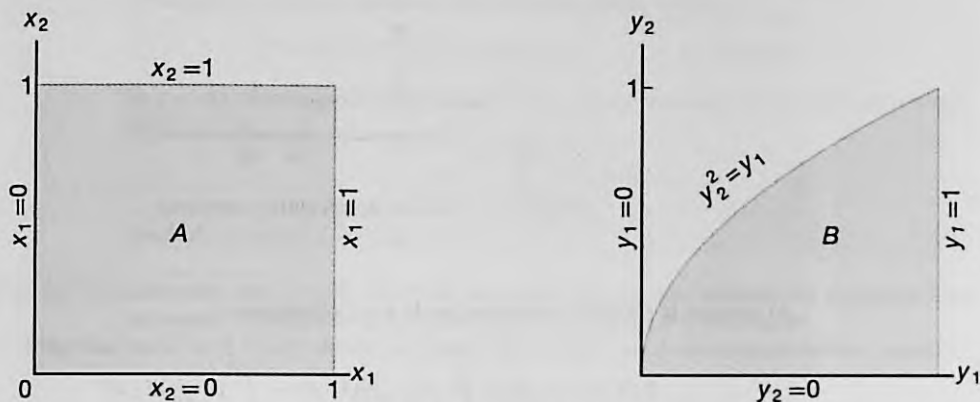


Figura 7.1: Gráfica del conjunto A en el conjunto B .

A menudo surgen problemas cuando deseamos encontrar la distribución de probabilidad de la variable aleatoria $Y = u(X)$ y X es una variable aleatoria continua y la transformación no es uno a uno. Es decir, a cada valor x le corresponde exactamente un valor y ; pero a cada valor y le corresponde más de un valor x . Por ejemplo, suponga que $f(x)$ es positiva en el intervalo $-1 < x < 2$ y cero en cualquier caso. Considere la transformación $y = x^2$. En este caso, $x = \pm \sqrt{y}$ para $0 < y < 1$ y $x = \sqrt{y}$ para $1 < y < 4$. Para el intervalo $1 < y < 4$, la distribución de probabilidad de Y se calcula como antes, con el teorema 7.3. Es decir,

$$g(y) = f[w(y)]|J| = \frac{f(\sqrt{y})}{2\sqrt{y}}, \quad 1 < y < 4.$$

Sin embargo, cuando $0 < y < 1$, podemos dividir el intervalo $-1 < x < 1$ para obtener las dos funciones inversas

$$x = -\sqrt{y}, \quad -1 < x < 0, \quad y \quad x = \sqrt{y}, \quad 0 < x < 1.$$

Entonces, a todo valor y le corresponde un solo valor x para cada partición. En la figura 7.2 vemos que

$$\begin{aligned} P(a < Y < b) &= P(-\sqrt{b} < X < -\sqrt{a}) + P(\sqrt{a} < X < \sqrt{b}) \\ &= \int_{-\sqrt{b}}^{-\sqrt{a}} f(x) dx + \int_{\sqrt{a}}^{\sqrt{b}} f(x) dx. \end{aligned}$$

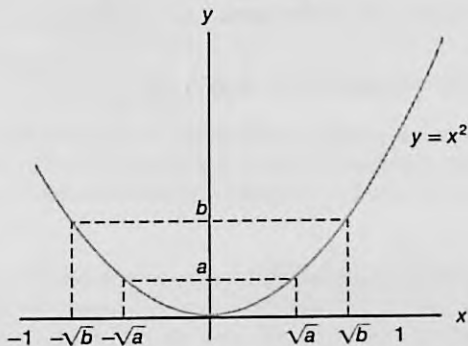


Figura 7.2: Función decreciente y creciente.

Al cambiar la variable de integración de x a y , obtenemos

$$\begin{aligned} P(a < Y < b) &= \int_b^a f(-\sqrt{y})J_1 dy + \int_a^b f(\sqrt{y})J_2 dy \\ &= - \int_a^b f(-\sqrt{y})J_1 dy + \int_a^b f(\sqrt{y})J_2 dy, \end{aligned}$$

donde

$$J_1 = \frac{d(-\sqrt{y})}{dy} = \frac{-1}{2\sqrt{y}} = -|J_1|$$

y

$$J_2 = \frac{d(\sqrt{y})}{dy} = \frac{1}{2\sqrt{y}} = |J_2|.$$

Por lo tanto, podremos escribir

$$P(a < Y < b) = \int_a^b [f(-\sqrt{y})|J_1| + f(\sqrt{y})|J_2|] dy,$$

y entonces

$$g(y) = f(-\sqrt{y})|J_1| + f(\sqrt{y})|J_2| = \frac{f(-\sqrt{y}) + f(\sqrt{y})}{2\sqrt{y}}, \quad 0 < y < 1.$$

La distribución de probabilidad de Y para $0 < y < 4$ se puede escribir ahora como

$$g(y) = \begin{cases} \frac{f(-\sqrt{y})+f(\sqrt{y})}{2\sqrt{y}}, & 0 < y < 1, \\ \frac{f(\sqrt{y})}{2\sqrt{y}}, & 1 < y < 4, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Este procedimiento para calcular $g(y)$ cuando $0 < y < 1$ se generaliza en el teorema 7.5 para k funciones inversas. Para transformaciones de funciones de diversas variables que no son uno a uno se recomienda al lector *Introduction to Mathematical Statistics* de Hogg, McKean y Craig (2005; véase la bibliografía).

Teorema 7.5: Suponga que X es una variable aleatoria continua con distribución de probabilidad $f(x)$. Definamos con $Y = u(X)$ una transformación entre los valores de X y Y que no es uno a uno. Si el intervalo sobre el que se define X se puede dividir en k conjuntos mutuamente disjuntos de manera que cada una de las funciones inversas

$$x_1 = w_1(y), \quad x_2 = w_2(y), \quad \dots, \quad x_k = w_k(y)$$

de $y = u(x)$ defina una correspondencia uno a uno, entonces la distribución de probabilidad de Y es

$$g(y) = \sum_{i=1}^k f[w_i(y)]|J_i|,$$

donde $J_i = w_i'(y)$, $i = 1, 2, \dots, k$.

Ejemplo 7.5: Demuestre que $Y = (X - \mu)^2/\sigma^2$ tiene una distribución chi cuadrada con 1 grado de libertad cuando X tiene una distribución normal con media μ y varianza σ^2 .

Solución: Sea $Z = (X - \mu)/\sigma$, donde la variable aleatoria Z tiene la distribución normal estándar

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, \quad -\infty < z < \infty.$$

Ahora debemos calcular la distribución de la variable aleatoria $Y = Z^2$. Las soluciones inversas de $y = z^2$ son $z = \pm\sqrt{y}$. Si designamos $z_1 = -\sqrt{y}$ y $z_2 = \sqrt{y}$, entonces $J_1 = -1/2\sqrt{y}$ y $J_2 = 1/2\sqrt{y}$. Entonces, por el teorema 7.5, tenemos

$$g(y) = \frac{1}{\sqrt{2\pi}} e^{-y/2} \left| \frac{-1}{2\sqrt{y}} \right| + \frac{1}{\sqrt{2\pi}} e^{-y/2} \left| \frac{1}{2\sqrt{y}} \right| = \frac{1}{\sqrt{2\pi}} y^{1/2-1} e^{-y/2}, \quad y > 0.$$

Como $g(y)$ es una función de densidad, se deduce que

$$1 = \frac{1}{\sqrt{2\pi}} \int_0^{\infty} y^{1/2-1} e^{-y/2} dy = \frac{\Gamma(1/2)}{\sqrt{\pi}} \int_0^{\infty} \frac{y^{1/2-1} e^{-y/2}}{\sqrt{2}\Gamma(1/2)} dy = \frac{\Gamma(1/2)}{\sqrt{\pi}},$$

la integral es el área bajo una curva de probabilidad gamma con los parámetros $\alpha = 1/2$ y $\beta = 2$. Por lo tanto, $\sqrt{\pi} = \Gamma(1/2)$ y la densidad de Y es dada por

$$g(y) = \begin{cases} \frac{1}{\sqrt{2}\Gamma(1/2)} y^{1/2-1} e^{-y/2}, & y > 0, \\ 0, & \text{en cualquier caso.} \end{cases}$$

que se considera una distribución chi cuadrada con 1 grado de libertad. ▮

7.3 Momentos y funciones generadoras de momentos

En esta sección nos concentramos en aplicaciones de las funciones generadoras de momentos. El propósito evidente de la función generadora de momentos es la determinación de los momentos de variables aleatorias. Sin embargo, la contribución más importante consiste en establecer distribuciones de funciones de variables aleatorias.

Si $g(X) = X^r$ para $r = 0, 1, 2, 3, \dots$, la definición 7.1 proporciona un valor esperado que se denomina r -ésimo momento alrededor del origen de la variable aleatoria X , que denotamos con μ'_r .

Definición 7.1: El r -ésimo momento alrededor del origen de la variable aleatoria X es dado por

$$\mu'_r = E(X^r) = \begin{cases} \sum x^r f(x), & \text{si } X \text{ es discreta,} \\ \int_{-\infty}^{\infty} x^r f(x) dx, & \text{si } X \text{ es continua.} \end{cases}$$

Como el primer y segundo momentos alrededor del origen son dados por $\mu'_1 = E(X)$ y $\mu'_2 = E(X^2)$, podemos escribir la media y la varianza de una variable aleatoria como

$$\mu = \mu'_1 \quad \text{y} \quad \sigma^2 = \mu'_2 - \mu^2.$$

Aunque los momentos de una variable aleatoria se pueden determinar directamente a partir de la definición 7.1, existe un procedimiento alternativo, el cual requiere que utilicemos una **función generadora de momentos**.

Definición 7.2: La **función generadora de momentos** de la variable aleatoria X es dada por $E(e^{tX})$, y se denota con $M_X(t)$. Por lo tanto,

$$M_X(t) = E(e^{tX}) = \begin{cases} \sum e^{tx} f(x), & \text{si } X \text{ es discreta,} \\ \int_{-\infty}^{\infty} e^{tx} f(x) dx, & \text{si } X \text{ es continua.} \end{cases}$$

Las funciones generadoras de momentos existirán sólo si la sumatoria o integral de la definición 7.2 converge. Si existe una función generadora de momentos de una variable aleatoria X , se puede utilizar para generar todos los momentos de dicha variable. El método se describe en el teorema 7.6 sin demostración.

Teorema 7.6: Sea X una variable aleatoria con función generadora de momentos $M_X(t)$. Entonces,

$$\left. \frac{d^r M_X(t)}{dt^r} \right|_{t=0} = \mu'_r.$$

Ejemplo 7.6: Calcule la función generadora de momentos de la variable aleatoria binomial X y después utilícela para verificar que $\mu = np$ y $\sigma^2 = npq$.

Solución: A partir de la definición 7.2 tenemos

$$M_X(t) = \sum_{x=0}^n e^{tx} \binom{n}{x} p^x q^{n-x} = \sum_{x=0}^n \binom{n}{x} (pe^t)^x q^{n-x}.$$

Al reconocer a esta última sumatoria como la expansión binomial de $(pe^t + q)^n$ obtenemos

$$M_X(t) = (pe^t + q)^n.$$

Así,

$$\frac{dM_X(t)}{dt} = n(pe^t + q)^{n-1}pe^t$$

y

$$\frac{d^2M_X(t)}{dt^2} = np[e^t(n-1)(pe^t + q)^{n-2}pe^t + (pe^t + q)^{n-1}e^t].$$

Al establecer $t = 0$ obtenemos

$$\mu'_1 = np \text{ y } \mu'_2 = np[(n-1)p + 1].$$

Por consiguiente,

$$\mu = \mu'_1 = np \text{ y } \sigma^2 = \mu'_2 - \mu^2 = np(1-p) = npq,$$

que coincide con los resultados que se obtuvieron en el capítulo 5. ▀

Ejemplo 7.7: Demuestre que la función generadora de momentos de la variable aleatoria X , la cual tiene una distribución de probabilidad normal con media μ y varianza σ^2 , es dada por

$$M_X(t) = \exp\left(\mu + \frac{1}{2}\sigma^2 t^2\right).$$

Solución: A partir de la definición 7.2, la función generadora de momentos de la variable aleatoria normal X es

$$\begin{aligned} M_X(t) &= \int_{-\infty}^{\infty} e^{tx} \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{x^2 - 2(\mu + t\sigma^2)x + \mu^2}{2\sigma^2}\right] dx. \end{aligned}$$

Si completamos el cuadrado en el exponente, podemos escribir

$$x^2 - 2(\mu + t\sigma^2)x + \mu^2 = [x - (\mu + t\sigma^2)]^2 - 2\mu t\sigma^2 - t^2\sigma^4$$

y, entonces,

$$\begin{aligned} M_X(t) &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{[x - (\mu + t\sigma^2)]^2 - 2\mu t\sigma^2 - t^2\sigma^4}{2\sigma^2}\right\} dx \\ &= \exp\left(\frac{2\mu t + \sigma^2 t^2}{2}\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{[x - (\mu + t\sigma^2)]^2}{2\sigma^2}\right\} dx. \end{aligned}$$

Sea $w = [x - (\mu + t\sigma^2)]/\sigma$; entonces $dx = \sigma dw$ y

$$M_X(t) = \exp\left(\mu + \frac{1}{2}\sigma^2 t^2\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-w^2/2} dw = \exp\left(\mu + \frac{1}{2}\sigma^2 t^2\right).$$

ya que la última integral representa el área bajo una curva de densidad normal estándar y, en consecuencia, es igual a 1.

Aunque el método de transformación de variables brinda una forma eficaz para determinar la distribución de una función de múltiples variables, existe un procedimiento alternativo, y que a menudo se prefiere cuando la función a analizar es una combinación lineal de variables aleatorias independientes. Este procedimiento utiliza las propiedades de las funciones generadoras de momentos que se estudian en los siguientes cuatro teoremas. Para no rebasar el alcance matemático de este libro, establecemos el teorema 7.7 sin demostración.

Teorema 7.7: (Teorema de unicidad) Sean X y Y dos variables aleatorias con funciones generadoras de momentos $M_X(t)$ y $M_Y(t)$, respectivamente. Si $M_X(t) = M_Y(t)$ para todos los valores de t , entonces X y Y tienen la misma distribución de probabilidad.

Teorema 7.8: $M_{X+a}(t) = e^{at} M_X(t)$.

Prueba: $M_{X+a}(t) = E[e^{t(X+a)}] = e^{at} E(e^{tX}) = e^{at} M_X(t)$.

Teorema 7.9: $M_{aX}(t) = M_X(at)$.

Prueba: $M_{aX}(t) = E[e^{t(aX)}] = E[e^{(at)X}] = M_X(at)$.

Teorema 7.10: Si X_1, X_2, \dots, X_n son variables aleatorias independientes con funciones generadoras de momentos $M_{X_1}(t), M_{X_2}(t), \dots, M_{X_n}(t)$, respectivamente, y $Y = X_1 + X_2 + \dots + X_n$, entonces,

$$M_Y(t) = M_{X_1}(t) M_{X_2}(t) \cdots M_{X_n}(t).$$

La demostración del teorema 7.10 se deja al lector.

Los teoremas 7.7 a 7.10 son fundamentales para entender las funciones generadoras de momentos. A continuación se presenta un ejemplo como ilustración. Hay muchas situaciones en que necesitamos conocer la distribución de la suma de las variables aleatorias. Podemos utilizar los teoremas 7.7 y 7.10, así como el resultado del ejercicio 7.19 de la página 224, para calcular la distribución de una suma de dos variables aleatorias independientes de Poisson, con funciones generadoras de momentos dadas por

$$M_{X_1}(t) = e^{\mu_1(e^t - 1)} \text{ y } M_{X_2}(t) = e^{\mu_2(e^t - 1)},$$

respectivamente. De acuerdo con el teorema 7.10, la función generadora de momentos de la variable aleatoria $Y_1 = X_1 + X_2$ es

$$M_{Y_1}(t) = M_{X_1}(t) M_{X_2}(t) = e^{\mu_1(e^t - 1)} e^{\mu_2(e^t - 1)} = e^{(\mu_1 + \mu_2)(e^t - 1)},$$

que de inmediato identificamos como la función generadora de momentos de una variable aleatoria que tiene una distribución de Poisson con el parámetro $\mu_1 + \mu_2$. Por lo tanto, de acuerdo con el teorema 7.7, de nuevo concluimos que la suma de dos variables aleatorias independientes, que tienen distribuciones de Poisson con los parámetros μ_1 y μ_2 , tiene una distribución de Poisson con el parámetro $\mu_1 + \mu_2$.

Combinaciones lineales de variables aleatorias

En estadística aplicada a menudo se necesita conocer la distribución de probabilidad de una combinación lineal de variables aleatorias normales independientes. Obtengamos la distribución de la variable aleatoria $Y = a_1X_1 + a_2X_2$ cuando X_1 es una variable normal con media μ_1 y varianza σ_1^2 y X_2 también es una variable normal, pero independiente de X_1 , con media μ_2 y varianza σ_2^2 . Primero, por medio del teorema 7.10, obtenemos

$$M_Y(t) = M_{a_1X_1}(t)M_{a_2X_2}(t),$$

y después, usando el teorema 7.9, obtenemos

$$M_Y(t) = M_{X_1}(a_1t)M_{X_2}(a_2t).$$

Si sustituimos a_1t por t , y después a_2t por t , en una función generadora de momentos de la distribución normal derivada en el ejemplo 7.7, tenemos

$$\begin{aligned} M_Y(t) &= \exp(a_1\mu_1t + a_1^2\sigma_1^2t^2/2 + a_2\mu_2t + a_2^2\sigma_2^2t^2/2) \\ &= \exp[(a_1\mu_1 + a_2\mu_2)t + (a_1^2\sigma_1^2 + a_2^2\sigma_2^2)t^2/2], \end{aligned}$$

que reconocemos como la función generadora de momentos de una distribución que es normal, con media $a_1\mu_1 + a_2\mu_2$ y varianza $a_1^2\sigma_1^2 + a_2^2\sigma_2^2$.

Al generalizar para el caso de n variables normales independientes, establecemos el siguiente resultado.

Teorema 7.11: Si X_1, X_2, \dots, X_n son variables aleatorias independientes que tienen distribuciones normales con medias $\mu_1, \mu_2, \dots, \mu_n$ y varianzas $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$, respectivamente, entonces la variable aleatoria

$$Y = a_1X_1 + a_2X_2 + \dots + a_nX_n$$

tiene una distribución normal con media

$$\mu_Y = a_1\mu_1 + a_2\mu_2 + \dots + a_n\mu_n$$

y varianza

$$\sigma_Y^2 = a_1^2\sigma_1^2 + a_2^2\sigma_2^2 + \dots + a_n^2\sigma_n^2.$$

Ahora es evidente que la distribución de Poisson y la distribución normal tienen una propiedad reproductiva, en el sentido de que la suma de variables aleatorias independientes que tengan cualquiera de estas distribuciones es una variable aleatoria que también tiene el mismo tipo de distribución. La distribución chi cuadrada también posee esta propiedad reproductiva.

Teorema 7.12: Si X_1, X_2, \dots, X_n son variables aleatorias mutuamente independientes, que tienen distribuciones chi cuadrada con v_1, v_2, \dots, v_n grados de libertad, respectivamente, entonces la variable aleatoria

$$Y = X_1 + X_2 + \dots + X_n$$

tiene una distribución chi cuadrada con $v = v_1 + v_2 + \dots + v_n$ grados de libertad.

Prueba: Por medio del teorema 7.10 y el ejercicio 7.21,

$$M_Y(t) = M_{X_1}(t)M_{X_2}(t) \cdots M_{X_n}(t) \text{ y } M_{X_i}(t) = (1 - 2t)^{-v_i/2}, \quad i = 1, 2, \dots, n.$$

Por lo tanto,

$$M_Y(t) = (1 - 2t)^{-v_1/2} (1 - 2t)^{-v_2/2} \dots (1 - 2t)^{-v_n/2} = (1 - 2t)^{-(v_1 + v_2 + \dots + v_n)/2},$$

que reconocemos como la función generadora de momentos de una distribución chi cuadrada con $v = v_1 + v_2 + \dots + v_n$ grados de libertad.

Corolario 7.1: Si X_1, X_2, \dots, X_n son variables aleatorias independientes que tienen distribuciones normales idénticas, con media μ y varianza σ^2 , entonces la variable aleatoria

$$Y = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2$$

tiene una distribución chi cuadrada con $v = n$ grados de libertad.

Este corolario es una consecuencia inmediata del ejemplo 7.5, y establece una relación entre la muy importante distribución chi cuadrada y la distribución normal. También debe brindar al lector una idea muy clara de lo que significa el parámetro llamado grados de libertad. En futuros capítulos el concepto de grados de libertad desempeñará un papel cada vez más relevante.

Corolario 7.2: Si X_1, X_2, \dots, X_n son variables aleatorias independientes y X_i tiene una distribución normal con media μ_i y varianza σ_i^2 para $i = 1, 2, \dots, n$, entonces la variable aleatoria

$$Y = \sum_{i=1}^n \left(\frac{X_i - \mu_i}{\sigma_i} \right)^2$$

tiene una distribución chi cuadrada con $v = n$ grados de libertad.

Ejercicios

7.1 Sea X una variable aleatoria que tiene la siguiente probabilidad

$$f(x) = \begin{cases} \frac{1}{3}, & x = 1, 2, 3, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la distribución de probabilidad de la variable aleatoria $Y = 2X - 1$.

7.2 Sea X una variable aleatoria binomial con la siguiente distribución de probabilidad

$$f(x) = \begin{cases} \binom{3}{x} \left(\frac{2}{5}\right)^x \left(\frac{3}{5}\right)^{3-x}, & x = 0, 1, 2, 3, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la distribución de probabilidad de la variable aleatoria $Y = X^2$.

7.3 Sean X_1 y X_2 variables aleatorias discretas con la siguiente distribución multinomial conjunta

$$f(x_1, x_2)$$

$$= \binom{2}{x_1, x_2, 2 - x_1 - x_2} \left(\frac{1}{4}\right)^{x_1} \left(\frac{1}{3}\right)^{x_2} \left(\frac{5}{12}\right)^{2 - x_1 - x_2}$$

para $x_1 = 0, 1, 2$; $x_2 = 0, 1, 2$; $x_1 + x_2 \leq 2$; y cero en cualquier caso. Calcule la distribución de probabilidad conjunta de $Y_1 = X_1 + X_2$ y $Y_2 = X_1 - X_2$.

7.4 Sean X_1 y X_2 variables aleatorias discretas con la siguiente distribución de probabilidad conjunta

$$f(x_1, x_2) = \begin{cases} \frac{3+3x_2}{18}, & x_1 = 1, 2; x_2 = 1, 2, 3, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la distribución de probabilidad de la variable aleatoria $Y = X_1 X_2$.

7.5 Si X tiene la siguiente distribución de probabilidad

$$f(x) = \begin{cases} 1, & 0 < x < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Demuestre que la variable aleatoria $Y = -2\ln X$ tiene una distribución chi cuadrada con 2 grados de libertad.

7.6 Dada la variable aleatoria X con la siguiente distribución de probabilidad

$$f(x) = \begin{cases} 2x, & 0 < x < 1, \\ 0, & \text{en cualquier caso,} \end{cases}$$

calcule la distribución de probabilidad de $Y = 8X^3$.

7.7 La velocidad de una molécula en un gas uniforme en equilibrio es una variable aleatoria V , cuya distribución de probabilidad es dada por

$$f(v) = \begin{cases} kv^2 e^{-bv^2}, & v > 0, \\ 0, & \text{en cualquier caso,} \end{cases}$$

donde k es una constante adecuada y b depende de la temperatura absoluta y de la masa de la molécula. Calcule la distribución de probabilidad de la energía cinética de la molécula W , donde $W = mV^2/2$.

7.8 La utilidad de un distribuidor, en unidades de \$5000, sobre un automóvil nuevo, es dada por $Y = X^2$, donde X es una variable aleatoria que tiene la siguiente función de densidad

$$f(x) = \begin{cases} 2(1-x), & 0 < x < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

- Calcule la función de densidad de probabilidad de la variable aleatoria Y .
- Utilice la función de densidad de Y para calcular la probabilidad de que la utilidad sobre el siguiente automóvil nuevo que venda este distribuidor sea menor que \$500.

7.9 El periodo hospitalario, en días, para pacientes que siguen un tratamiento para cierto tipo de enfermedad del riñón es una variable aleatoria $Y = X + 4$, donde X tiene la siguiente función de densidad

$$f(x) = \begin{cases} \frac{32}{(x+4)^3}, & x > 0, \\ 0, & \text{en cualquier caso.} \end{cases}$$

- Calcule la función de densidad de probabilidad de la variable aleatoria Y .
- Utilice la función de densidad de Y para calcular la probabilidad de que el periodo hospitalario para un paciente que sigue este tratamiento exceda los 8 días.

7.10 Las variables aleatorias X y Y , que representan los pesos de cremas y chiclosos, respectivamente, en

cajas de un kilogramo de chocolates que contienen una combinación de cremas, chiclosos y envinados, tienen la siguiente función de densidad conjunta

$$f(x, y) = \begin{cases} 24xy, & 0 \leq x \leq 1, 0 \leq y \leq 1, x + y \leq 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

- Calcule la función de densidad de probabilidad de la variable aleatoria $Z = X + Y$.
- Utilice la función de densidad de Z para calcular la probabilidad de que, en una determinada caja, la suma de los pesos de las cremas y los chiclosos sea por lo menos $1/2$ del peso total, pero menos de $3/4$.

7.11 La cantidad de queroseno en un tanque al inicio de cualquier día, en miles de litros, es una cantidad aleatoria Y , de la cual una cantidad aleatoria X se vende durante ese día. Suponga que la función de densidad conjunta de estas variables es dada por

$$f(x, y) = \begin{cases} 2, & 0 < x < y, 0 < y < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la función de densidad de probabilidad para la cantidad de queroseno que queda en el tanque al final del día.

7.12 Sean X_1 y X_2 variables aleatorias independientes que tienen cada una la siguiente distribución de probabilidad

$$f(x) = \begin{cases} e^{-x}, & x > 0, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Demuestre que las variables aleatorias Y_1 y Y_2 son independientes cuando $Y_1 = X_1 + X_2$ y $Y_2 = X_1/(X_1 + X_2)$.

7.13 Una corriente de I amperios que fluye a través de una resistencia de R ohms varía de acuerdo con la siguiente distribución de probabilidad

$$f(i) = \begin{cases} 6i(1-i), & 0 < i < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Si la resistencia varía independientemente de la corriente de acuerdo con la siguiente distribución de probabilidad

$$g(r) = \begin{cases} 2r, & 0 < r < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

calcule la distribución de probabilidad para la potencia $W = I^2R$ watts.

7.14 Sea X una variable aleatoria con la siguiente distribución de probabilidad

$$f(x) = \begin{cases} \frac{1+x}{2}, & -1 < x < 1, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la distribución de probabilidad de la variable aleatoria $Y = X^2$.

7.15 Si X tiene la siguiente distribución de probabilidad

$$f(x) = \begin{cases} \frac{2(x+1)}{9}, & -1 < x < 2, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Calcule la distribución de probabilidad de la variable aleatoria $Y = X^2$.

7.16 Demuestre que el r -ésimo momento respecto al origen de la distribución gamma es

$$\mu'_r = \frac{\beta^r \Gamma(\alpha + r)}{\Gamma(\alpha)}.$$

[Sugerencia: Sustituya $y = x/\beta$ en la integral que define μ'_r y después utilice la función gamma para evaluar la integral].

7.17 Una variable aleatoria X tiene la siguiente distribución uniforme discreta

$$f(x; k) = \begin{cases} \frac{1}{k}, & x = 1, 2, \dots, k, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Demuestre que la función generadora de momentos de X es

$$M_X(t) = \frac{e^t(1 - e^{kt})}{k(1 - e^t)}.$$

7.18 Una variable aleatoria X tiene la distribución geométrica $g(x; p) = pq^{x-1}$ para $x = 1, 2, 3, \dots$. Demuestre que la función generadora de momentos de X es

$$M_X(t) = \frac{pe^t}{1 - qe^t}, \quad t < \ln q,$$

y después use $M_X(t)$ para calcular la media y la varianza de la distribución geométrica.

7.19 Una variable aleatoria X tiene la distribución de Poisson $p(x; \mu) = e^{-\mu} \mu^x / x!$ para $x = 0, 1, 2, \dots$. Demuestre que la función generadora de momentos de X es

$$M_X(t) = e^{\mu(e^t - 1)}.$$

Utilice $M_X(t)$ para calcular la media y la varianza de la distribución de Poisson.

7.20 La función generadora de momentos de cierta variable aleatoria de Poisson X es dada por

$$M_X(t) = e^{4(e^t - 1)}.$$

Calcule $P(\mu - 2\sigma < X < \mu + 2\sigma)$.

7.21 Demuestre que la función generadora de momentos de la variable aleatoria X , que tiene una distribución chi cuadrada con ν grados de libertad, es

$$M_X(t) = (1 - 2t)^{-\nu/2}.$$

7.22 Con la función generadora de momentos del ejemplo 7.21 demuestre que la media y la varianza de la distribución chi cuadrada con ν grados de libertad son, respectivamente, ν y 2ν .

7.23 Si tanto X como Y , distribuidas de manera independiente, siguen distribuciones exponenciales con parámetro medio 1, calcule las distribuciones de

- a) $U = X + Y$;
b) $V = X/(X + Y)$.

7.24 Mediante la expansión de e^{tx} en una serie de Maclaurin y la integración término por término, demuestre que

$$\begin{aligned} M_X(t) &= \int_{-\infty}^{\infty} e^{tx} f(x) dx \\ &= 1 + \mu t + \mu'_2 \frac{t^2}{2!} + \dots + \mu'_r \frac{t^r}{r!} + \dots \end{aligned}$$

Capítulo 8

Distribuciones de muestreo fundamentales y descripciones de datos

8.1 Muestreo aleatorio

El resultado de un experimento estadístico se puede registrar como un valor numérico o como una representación descriptiva. Cuando se lanza un par de dados y lo que nos interesa es el resultado total, registramos un valor numérico. Sin embargo, si a los estudiantes de cierta escuela se les hacen pruebas de sangre para averiguar cuál es su tipo, podría ser más útil una representación descriptiva. La sangre de una persona se puede clasificar de 8 maneras. Puede ser AB, A, B u O, cada una con un signo de más o de menos, lo cual depende de la presencia o ausencia del antígeno Rh.

En este capítulo nos enfocamos en el muestreo de distribuciones o poblaciones, y estudiamos cantidades tan importantes como la *media de la muestra* y la *varianza de la muestra*, que serán de importancia fundamental en los capítulos siguientes. Además, en los próximos capítulos intentamos introducir al lector al papel que desempeñarán la media y la varianza de la muestra en la inferencia estadística. El uso de las computadoras modernas de alta velocidad permite a los científicos e ingenieros incrementar enormemente su uso de la inferencia estadística formal con técnicas gráficas. La mayoría de las veces la inferencia formal parece muy árida y quizás incluso abstracta para el profesional o el gerente que desea que el análisis estadístico sea una guía para la toma de decisiones.

Poblaciones y muestras

Comenzamos esta sección presentando los conceptos de *poblaciones* y *muestras*. Ambas se mencionan de forma extensa en el capítulo 1; sin embargo, aquí será necesario estudiarlas más ampliamente, en particular en el contexto del concepto de variables aleatorias. La totalidad de observaciones que nos interesan, ya sean de número finito o infinito, constituye lo que llamamos **población**. En alguna época el término *población* se refería a observaciones que se obtenían de estudios estadísticos aplicados a personas. En la actualidad el estadístico utiliza la palabra para referirse a observaciones sobre cualquier cuestión de interés, ya sea de grupos de personas, de animales o de todos los resultados posibles de algún complicado sistema biológico o de ingeniería.

Definición 8.1: Una población consta de la totalidad de las observaciones en las que estamos interesados.

El número de observaciones en la población se define como el tamaño de la población. Si en la escuela hay 600 estudiantes que clasificamos de acuerdo con su tipo de sangre, decimos que tenemos una población de tamaño 600. Los números en las cartas de una baraja, las estaturas de los residentes de cierta ciudad y las longitudes de los peces en un lago específico son ejemplos de poblaciones de tamaño finito. En cada caso el número total de observaciones es un número finito. Las observaciones que se obtienen al medir diariamente la presión atmosférica desde el pasado hasta el futuro, o todas las mediciones de la profundidad de un lago desde cualquier posición concebible son ejemplos de poblaciones cuyos tamaños son infinitos. Algunas poblaciones finitas son tan grandes que en teoría las supondríamos infinitas, lo cual es cierto si se considera la población de la vida útil de cierto tipo de batería de almacenamiento que se está fabricando para distribuirla en forma masiva en todo el país.

Cada observación en una población es un valor de una variable aleatoria X que tiene alguna distribución de probabilidad $f(x)$. Si se inspeccionan artículos que salen de una línea de ensamble para buscar defectos, entonces cada observación en la población podría ser un valor 0 o 1 de la variable aleatoria X de Bernoulli, con una distribución de probabilidad

$$b(x; 1, p) = p^x q^{1-x}, \quad x = 0, 1$$

donde 0 indica un artículo sin defecto y 1 indica un artículo defectuoso. De hecho, se supone que p , la probabilidad de que cualquier artículo esté defectuoso, permanece constante de una prueba a otra. En el experimento del tipo de sangre la variable aleatoria X representa el tipo de sangre y se supone que toma un valor del 1 al 8. A cada estudiante se le asigna uno de los valores de la variable aleatoria discreta. Las duraciones de las baterías de almacenamiento son valores que toma una variable aleatoria continua que quizá tiene una distribución normal. De ahora en adelante, cuando nos refiramos a una "población binomial", a una "población normal" o, en general, a la "población $f(x)$ ", aludiremos a una población cuyas observaciones son valores de una variable aleatoria que tiene una distribución binomial, una distribución normal o la distribución de probabilidad $f(x)$. Por ello, a la media y a la varianza de una variable aleatoria o distribución de probabilidad también se les denomina la media y la varianza de la población correspondiente.

En el campo de la inferencia estadística, el estadístico se interesa en llegar a conclusiones respecto a una población, cuando es imposible o poco práctico conocer todo el conjunto de observaciones que la constituyen. Por ejemplo, al intentar determinar la longitud de la vida promedio de cierta marca de bombilla, sería imposible probarlas todas si tenemos que dejar algunas para venderlas. Los costos desmesurados que implicaría estudiar a toda la población también constituirían un factor que impediría hacerlo. Por lo tanto, debemos depender de un subconjunto de observaciones de la población que nos ayude a realizar inferencias respecto a ella. Esto nos lleva a considerar el concepto de muestreo.

Definición 8.2: Una muestra es un subconjunto de una población.

Para que las inferencias que hacemos sobre la población a partir de la muestra sean válidas, debemos obtener muestras que sean representativas de ella. Con mucha

frecuencia nos sentimos tentados a elegir una muestra seleccionando a los miembros más convenientes de la población. Tal procedimiento podría conducir a inferencias erróneas respecto a la población. Se dice que cualquier procedimiento de muestreo que produzca inferencias que sobreestimen o subestimen de forma consistente alguna característica de la población está **sesgado**. Para eliminar cualquier posibilidad de sesgo en el procedimiento de muestreo es deseable elegir una **muestra aleatoria**, lo cual significa que las observaciones se realicen de forma independiente y al azar.

Para seleccionar una muestra aleatoria de tamaño n de una población $f(x)$ definimos la variable aleatoria X_i , $i = 1, 2, \dots, n$, que representa la i -ésima medición o valor de la muestra que observamos. Si las mediciones se obtienen repitiendo el experimento n veces independientes en, esencialmente, las mismas condiciones, las variables aleatorias X_1, X_2, \dots, X_n constituirán entonces una muestra aleatoria de la población $f(x)$ con valores numéricos x_1, x_2, \dots, x_n . Debido a las condiciones idénticas en las que se seleccionan los elementos de la muestra, es razonable suponer que las n variables aleatorias X_1, X_2, \dots, X_n son independientes y que cada una tiene la misma distribución de probabilidad $f(x)$. Es decir, las distribuciones de probabilidad de X_1, X_2, \dots, X_n son, respectivamente, $f(x_1), f(x_2), \dots, f(x_n)$, y su distribución de probabilidad conjunta es $f(x_1, x_2, \dots, x_n) = f(x_1) f(x_2) \cdots f(x_n)$. El concepto de muestra aleatoria se describe de manera formal en la siguiente definición.

Definición 8.3: Sean X_1, X_2, \dots, X_n variables aleatorias independientes n , cada una con la misma distribución de probabilidad $f(x)$. Definimos X_1, X_2, \dots, X_n como una **muestra aleatoria** de tamaño n de la población $f(x)$ y escribimos su distribución de probabilidad conjunta como

$$f(x_1, x_2, \dots, x_n) = f(x_1) f(x_2) \cdots f(x_n).$$

Si se realiza una selección aleatoria de $n = 8$ baterías de almacenamiento de un proceso de fabricación que mantiene las mismas especificaciones, y al registrar la duración de cada batería se encuentra que la primera medición x_1 es un valor de X_1 , la segunda medición x_2 es un valor de X_2 , y así sucesivamente, entonces x_1, x_2, \dots, x_8 son los valores de la muestra aleatoria X_1, X_2, \dots, X_8 . Si suponemos que la población de vidas útiles de las baterías es normal, los valores posibles de cualquier X_i , $i = 1, 2, \dots, 8$ serán exactamente los mismos que los de la población original, por consiguiente, X_i tiene una distribución normal idéntica a la de X .

8.2 Algunos estadísticos importantes

Nuestro principal propósito al seleccionar muestras aleatorias consiste en obtener información acerca de los parámetros desconocidos de la población. Suponga, por ejemplo, que deseamos concluir algo respecto a la proporción de consumidores de café en Estados Unidos que prefieren cierta marca de café. Sería imposible interrogar a cada consumidor estadounidense de café para calcular el valor del parámetro p que representa la proporción de la población. En vez de esto se selecciona una muestra aleatoria grande y se calcula la proporción \hat{p} de personas en esta muestra que prefieren la marca de café en cuestión. El valor \hat{p} se utiliza ahora para hacer una inferencia respecto a la proporción p verdadera.

Ahora, \hat{p} es una función de los valores observados en la muestra aleatoria; ya que es posible tomar muchas muestras aleatorias de la misma población, esperaríamos

que \hat{p} variara un poco de una a otra muestra. Es decir, \hat{p} es un valor de una variable aleatoria que representamos con P . Tal variable aleatoria se llama **estadístico**.

Definición 8.4: Cualquier función de las variables aleatorias que forman una muestra aleatoria se llama **estadístico**.

Medidas de localización de una muestra: la media, la mediana y la moda muestrales

En el capítulo 4 presentamos los parámetros μ y σ^2 , que miden el centro y la variabilidad de una distribución de probabilidad. Éstos son parámetros de población constantes y de ninguna manera se ven afectados o influidos por las observaciones de una muestra aleatoria. Definiremos, sin embargo, algunos estadísticos importantes que describen las medidas correspondientes de una muestra aleatoria. Los estadísticos que más se utilizan para medir el centro de un conjunto de datos, acomodados en orden de magnitud, son la **media**, la **mediana** y la **moda**. Aunque los primeros dos estadísticos se expusieron en el capítulo 1, repetiremos las definiciones. Sean X_1, X_2, \dots, X_n representaciones de n variables aleatorias.

a) Media muestral:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Observe que el estadístico \bar{X} toma el valor $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ cuando X_1 toma el valor x_1 , X_2 toma el valor x_2 y así sucesivamente. El término *media muestral* se aplica tanto al estadístico \bar{X} como a su valor calculado \bar{x} .

b) Mediana muestral:

$$\bar{x} = \begin{cases} x_{(n+1)/2}, & \text{si } n \text{ es impar,} \\ \frac{1}{2}(x_{n/2} + x_{n/2+1}), & \text{si } n \text{ es par.} \end{cases}$$

La mediana muestral también es una medida de localización que indica el valor central de la muestra. En la sección 1.3 se presentan ejemplos de la media muestral y de la mediana muestral. La moda muestral se define de la siguiente manera:

c) La moda muestral es el valor que ocurre con mayor frecuencia en la muestra.

Ejemplo 8.1: Suponga que un conjunto de datos consta de las siguientes observaciones:

0.32 0.53 0.28 0.37 0.47 0.43 0.36 0.42 0.38 0.43

La moda de la muestra es 0.43, ya que este valor aparece con más frecuencia que los demás.

Como se expuso en el capítulo 1, una medida de localización o tendencia central en una muestra no da por sí misma una indicación clara de la naturaleza de ésta, de manera que también debe considerarse una medida de variabilidad en la muestra.

Las medidas de variabilidad de una muestra: la varianza, la desviación estándar y el rango de la muestra

La variabilidad en la muestra refleja cómo se dispersan las observaciones a partir del promedio. Se remite al lector al capítulo 1 para un análisis más amplio. Es posible tener dos conjuntos de observaciones con las mismas media o mediana que difieran de manera considerable en la variabilidad de sus mediciones sobre el promedio.

Considere las siguientes mediciones, en litros, para dos muestras de jugo de naranja envasado por las empresas A y B:

Muestra A	0.97	1.00	0.94	1.03	1.06
Muestra B	1.06	1.01	0.88	0.91	1.14

Ambas muestras tienen la misma media, 1.00 litros. Es muy evidente que la empresa A envasa el jugo de naranja con un contenido más uniforme que la B. Decimos que la **variabilidad** o la **dispersión** de las observaciones a partir del promedio es menor para la muestra A que para la muestra B. Por lo tanto, al comprar jugo de naranja, tendríamos más confianza en que el envase que seleccionemos se acerque al promedio anunciado si se lo compramos a la empresa A.

En el capítulo 1 presentamos varias medidas de la variabilidad de una muestra, como la **varianza muestral**, la **desviación estándar muestral** y el **rango de la muestra**. En este capítulo nos enfocaremos sobre todo en la varianza de la muestra. Nuevamente, sea que X_1, X_2, \dots, X_n representan n variables aleatorias.

a) La varianza muestral:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2. \quad (8.2.1)$$

El valor calculado de S^2 para una muestra dada se denota con s^2 . Observe que S^2 se define esencialmente como el promedio de los cuadrados de las desviaciones de las observaciones a partir de su media. La razón para utilizar $n-1$ como divisor, en vez de la elección más obvia n , quedará más clara en el capítulo 9.

Ejemplo 8.2: Una comparación de los precios de café en 4 tiendas de abarrotes de San Diego, seleccionadas al azar, mostró aumentos en comparación con el mes anterior de 12, 15, 17 y 20 centavos por bolsa de una libra. Calcule la varianza de esta muestra aleatoria de aumentos de precio.

Solución: Si calculamos la media de la muestra, obtenemos

$$\bar{x} = \frac{12 + 15 + 17 + 20}{4} = 16 \text{ centavos.}$$

Por lo tanto,

$$\begin{aligned} s^2 &= \frac{1}{3} \sum_{i=1}^4 (x_i - 16)^2 = \frac{(12 - 16)^2 + (15 - 16)^2 + (17 - 16)^2 + (20 - 16)^2}{3} \\ &= \frac{(-4)^2 + (-1)^2 + (1)^2 + (4)^2}{3} = \frac{34}{3}. \end{aligned}$$

Mientras que la expresión para la varianza de la muestra de la definición 8.6 ilustra mejor que S^2 es una medida de variabilidad, una expresión alternativa tiene cierto mérito, de manera que el lector debería conocerla. El siguiente teorema contiene tal expresión.

Teorema 8.1: Si S^2 es la varianza de una muestra aleatoria de tamaño n , podemos escribir

$$S^2 = \frac{1}{n(n-1)} \left[n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2 \right].$$

Prueba: Por definición,

$$\begin{aligned} S^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i^2 - 2\bar{X}X_i + \bar{X}^2) \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2 \right]. \end{aligned}$$

Como en el capítulo 1, a continuación se definen la **desviación estándar muestral** y el **rango muestral**:

b) Desviación estándar muestral:

$$S = \sqrt{S^2},$$

donde S^2 es la varianza muestral.

Permitamos que X_{\max} denote el más grande de los valores X_i y X_{\min} el más pequeño.

c) Rango muestral:

$$R = X_{\max} - X_{\min}.$$

Ejemplo 8.3: Calcule la varianza de los datos 3, 4, 5, 6, 6 y 7, que representan el número de truchas atrapadas por una muestra aleatoria de 6 pescadores, el 19 de junio de 1996, en el lago Muskoka.

Solución: Encontramos que $\sum_{i=1}^6 x_i^2 = 171$, $\sum_{i=1}^6 x_i = 31$ y $n = 6$. De aquí,

$$s^2 = \frac{1}{(6)(5)} [(6)(171) - (31)^2] = \frac{13}{6}.$$

Por consiguiente, la desviación estándar de la muestra $s = \sqrt{13/6} = 1.47$ y el rango muestral es $7 - 3 = 4$.

Ejercicios

8.1 Defina las poblaciones adecuadas a partir de las cuales se seleccionaron las siguientes muestras:

- Se llamó por teléfono a personas de 200 casas en la ciudad de Richmond y se les pidió nombrar al candidato por el que votarían en la elección del presidente de la mesa directiva de la escuela.
- Se lanzó 100 veces una moneda y se registraron 34 cruces.

- Se probaron 200 pares de un nuevo tipo de calzado deportivo en un torneo de tenis profesional para determinar su duración y se encontró que, en promedio, duraron 4 meses.
- En cinco ocasiones diferentes a una abogada le tomó 21, 26, 24, 22 y 21 minutos conducir desde su casa en los suburbios hasta su oficina en el centro de la ciudad.

8.2 El tiempo, en minutos, que 10 pacientes esperan en un consultorio médico antes de recibir tratamiento se registraron como sigue: 5, 11, 9, 5, 10, 15, 6, 10, 5 y 10. Trate los datos como una muestra aleatoria y calcule

- la media;
- la mediana;
- la moda.

8.3 Los tiempos que los 9 individuos de una muestra aleatoria tardan en reaccionar ante un estimulante se registraron como 2.5, 3.6, 3.1, 4.3, 2.9, 2.3, 2.6, 4.1 y 3.4 segundos. Calcule

- la media;
- la mediana.

8.4 El número de multas emitidas por infracciones de tránsito por 8 oficiales estatales durante el fin de semana del día en Conmemoración de los Caídos es 5, 4, 7, 7, 6, 3, 8 y 6.

- Si estos valores representan el número de multas emitidas por una muestra aleatoria de 8 oficiales estatales del condado de Montgomery, en Virginia, defina una población adecuada.
- Si los valores representan el número de multas emitidas por una muestra aleatoria de 8 oficiales estatales de Carolina del Sur, defina una población adecuada.

8.5 El número de respuestas incorrectas en un examen de competencia de verdadero-falso para una muestra aleatoria de 15 estudiantes se registraron de la siguiente manera: 2, 1, 3, 0, 1, 3, 6, 0, 3, 3, 5, 2, 1, 4 y 2. Calcule

- la media;
- la mediana;
- la moda.

8.6 Calcule la media, la mediana y la moda para la muestra, cuyas observaciones, 15, 7, 8, 95, 19, 12, 8, 22 y 14 representan el número de días de incapacidad médica reportados en 9 solicitudes de devolución de impuestos. ¿Qué valor parece ser la mejor medida del centro de esos datos? Explique las razones de su preferencia.

8.7 Una muestra aleatoria de empleados de una fábrica local prometieron los siguientes donativos, en dólares, al United Fund: 100, 40, 75, 15, 20, 100, 75, 50, 30, 10, 55, 75, 25, 50, 90, 80, 15, 25, 45 y 100. Calcule

- la media;
- la moda.

8.8 De acuerdo con la escritora ecologista Jacqueline Killeen, los fosfatos que contienen los detergentes de uso casero pasan directamente a nuestros sistemas de desagüe, ocasionando que los lagos se conviertan

en pantanos, los cuales a la larga se volverán desiertos. Los siguientes datos muestran la cantidad de fosfatos por carga de lavado, en gramos, para una muestra aleatoria de diversos tipos de detergentes que se usan de acuerdo con las instrucciones prescritas:

Detergente para ropa	Fosfatos por carga (gramos)
A & P Blue Sail	48
Dash	47
Concentrated All	42
Cold Water All	42
Breeze	41
Oxydol	34
Ajax	31
Sears	30
Fab	29
Cold Power	29
Bold	29
Rinso	26

Para los datos de fosfato dados, calcule

- la media;
- la mediana;
- la moda.

8.9 Considere los datos del ejercicio 8.2 y calcule

- el rango;
- la desviación estándar.

8.10 Para la muestra de tiempos de reacción del ejercicio 8.3 calcule

- el rango;
- la varianza, utilizando la fórmula de la forma (8.2.1).

8.11 Para los datos del ejercicio 8.5 calcule la varianza utilizando la fórmula

- de la forma (8.2.1);
- del teorema 8.1.

8.12 El contenido de alquitrán de 8 marcas de cigarrillos que se seleccionan al azar de la lista más reciente publicada por la Comisión Federal de Comercio es el siguiente: 7.3, 8.6, 10.4, 16.1, 12.2, 15.1, 14.5 y 9.3 miligramos. Calcule

- la media;
- la varianza.

8.13 Los promedios de calificaciones de 20 estudiantes universitarios del último año, seleccionados al azar de una clase que se va a graduar, son los siguientes:

3.2	1.9	2.7	2.4	2.8
2.9	3.8	3.0	2.5	3.3
1.8	2.5	3.7	2.8	2.0
3.2	2.3	2.1	2.5	1.9

Calcule la desviación estándar.

8.14 a) Demuestre que la varianza de la muestra permanece sin cambio si a cada valor de la muestra se le suma o se le resta una constante c .

b) Demuestre que la varianza de la muestra se vuelve c^2 veces su valor original si cada observación de la muestra se multiplica por c .

8.15 Verifique que la varianza de la muestra 4, 9, 3, 6, 4 y 7 es 5.1, y utilice este hecho, junto con los resultados del ejercicio 8.14, para calcular

- a) la varianza de la muestra 12, 27, 9, 18, 12 y 21;
b) la varianza de la muestra 9, 14, 8, 11, 9 y 12.

8.16 En la temporada 2004-2005 el equipo de fútbol americano de la Universidad del Sur de California tuvo las siguientes diferencias de puntuación en los 13 partidos que jugó.

11 49 32 3 6 38 38 30 8 4 31 5 36

Calcule

- a) la media de la diferencia de puntos;
b) la mediana de las diferencias de puntos.

8.3 Distribuciones muestrales

El campo de la inferencia estadística trata básicamente con generalizaciones y predicciones. Por ejemplo, con base en las opiniones de varias personas entrevistadas en la calle, los estadounidenses podrían afirmar que en una próxima elección 60% de los votantes de la ciudad de Detroit favorecerían a cierto candidato. En este caso tratamos con una muestra aleatoria de opiniones de una población finita muy grande. Por otro lado, con base en las estimaciones de 3 contratistas seleccionados al azar, de los 30 que laboran actualmente en esta ciudad, podríamos afirmar que el costo promedio de construir una residencia en Charleston, Carolina del Sur, está entre \$330,000 y \$335,000. La población que se va a muestrear aquí también es finita, pero muy pequeña. Finalmente, consideremos una máquina despachadora de bebida gaseosa que está diseñada para servir en promedio 240 mililitros de bebida. Un ejecutivo de la empresa calcula la media de 40 bebidas servidas y obtiene $\bar{x} = 236$ mililitros y, con base en este valor, decide que la máquina está sirviendo bebidas con un contenido promedio de $\mu = 240$ mililitros. Las 40 bebidas servidas representan una muestra de la población infinita de posibles bebidas que despachará esta máquina.

Inferencias sobre la población a partir de información de la muestra

En cada uno de los ejemplos anteriores calculamos un estadístico de una muestra que se selecciona de la población, y con base en tales estadísticos hicimos varias afirmaciones respecto a los valores de los parámetros de la población, que pueden ser o no ciertas. El ejecutivo de la empresa decide que la máquina despachadora está sirviendo bebidas con un contenido promedio de 240 mililitros, aunque la media de la muestra fue de 236 mililitros, porque conoce la teoría del muestreo según la cual, si $\mu = 240$ mililitros, tal valor de la muestra podría ocurrir fácilmente. De hecho, si realiza pruebas similares, cada hora por ejemplo, esperaríamos que los valores del estadístico \bar{x} fluctuaran por arriba y por abajo de $\mu = 240$ mililitros. Sólo cuando el valor de \bar{x} difiera considerablemente de 240 mililitros el ejecutivo de la empresa tomará medidas para ajustar la máquina.

Como un estadístico es una variable aleatoria que depende sólo de la muestra observada, debe tener una distribución de probabilidad.

Definición 8.5: La distribución de probabilidad de un estadístico se denomina **distribución muestral**.

La distribución muestral de un estadístico depende de la distribución de la población, del tamaño de las muestras y del método de selección de las muestras. En lo que resta de este capítulo estudiaremos varias de las distribuciones muestrales más importantes de los estadísticos que se utilizan con frecuencia. Las aplicaciones de tales distribuciones muestrales a problemas de inferencia estadística se consideran en la mayoría de los capítulos posteriores. La distribución de probabilidad de \bar{X} se llama **distribución muestral de la media**.

¿Qué es la distribución muestral de \bar{X} ?

Se deberían considerar las distribuciones muestrales de \bar{X} y S^2 como los mecanismos a partir de los cuales se puede hacer inferencias acerca de los parámetros μ y σ^2 . La distribución muestral de \bar{X} con tamaño muestral n es la distribución que resulta cuando un experimento se lleva a cabo una y otra vez (siempre con una muestra de tamaño n) y resultan los diversos valores de \bar{X} . Por lo tanto, esta distribución muestral describe la variabilidad de los promedios muestrales alrededor de la media de la población μ . En el caso de la máquina despachadora de bebidas, el conocer la distribución muestral de \bar{X} le permite al analista encontrar una discrepancia "típica" entre un valor \bar{x} observado y el verdadero valor de μ . Se aplica el mismo principio en el caso de la distribución de S^2 . La distribución muestral produce información acerca de la variabilidad de los valores de s^2 alrededor de σ^2 en experimentos que se repiten.

8.4 Distribución muestral de medias y el teorema del límite central

La primera distribución muestral importante a considerar es la de la media \bar{X} . Suponga que de una población normal con media μ y varianza σ^2 se toma una muestra aleatoria de n observaciones. Cada observación X_i , $i = 1, 2, \dots, n$, de la muestra aleatoria tendrá entonces la misma distribución normal que la población de donde se tomó. Así, por la propiedad reproductiva de la distribución normal que se estableció en el teorema 7.11, concluimos que

$$\bar{X} = \frac{1}{n}(X_1 + X_2 + \dots + X_n)$$

tiene una distribución normal con media

$$\mu_{\bar{X}} = \frac{1}{n}(\underbrace{\mu + \mu + \dots + \mu}_{n \text{ términos}}) = \mu \text{ y varianza } \sigma_{\bar{X}}^2 = \frac{1}{n^2}(\underbrace{\sigma^2 + \sigma^2 + \dots + \sigma^2}_{n \text{ términos}}) = \frac{\sigma^2}{n}.$$

Si tomamos muestras de una población con distribución desconocida, ya sea finita o infinita, la distribución muestral de \bar{X} aún será aproximadamente normal con media μ y varianza σ^2/n , siempre que el tamaño de la muestra sea grande. Este asombroso resultado es una consecuencia inmediata del siguiente teorema, que se conoce como teorema del límite central.

El teorema del límite central

Teorema 8.2: **Teorema del límite central:** Si \bar{X} es la media de una muestra aleatoria de tamaño n , tomada de una población con media μ y varianza finita σ^2 , entonces la forma límite de la distribución de

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}},$$

a medida que $n \rightarrow \infty$, es la distribución normal estándar $n(z; 0, 1)$.

La aproximación normal para \bar{X} por lo general será buena si $n \geq 30$, siempre y cuando la distribución de la población no sea muy asimétrica. Si $n < 30$, la aproximación será buena sólo si la población no es muy diferente de una distribución normal y, como antes se estableció, si se sabe que la población es normal, la distribución muestral de \bar{X} seguirá siendo una distribución normal exacta, sin importar qué tan pequeño sea el tamaño de las muestras.

El tamaño de la muestra $n = 30$ es un lineamiento para el teorema del límite central. Sin embargo, como indica el planteamiento del teorema, la suposición de normalidad en la distribución de \bar{X} se vuelve más precisa a medida que n se hace más grande. De hecho, la figura 8.1 ilustra cómo funciona el teorema. La figura indica cómo la distribución de \bar{X} se acerca más a la normalidad a medida que aumenta n , empezando con la distribución claramente asimétrica de una observación individual ($n = 1$). También ilustra que la media de \bar{X} sigue siendo μ para cualquier tamaño de la muestra y que la varianza de \bar{X} se vuelve más pequeña a medida que aumenta n .

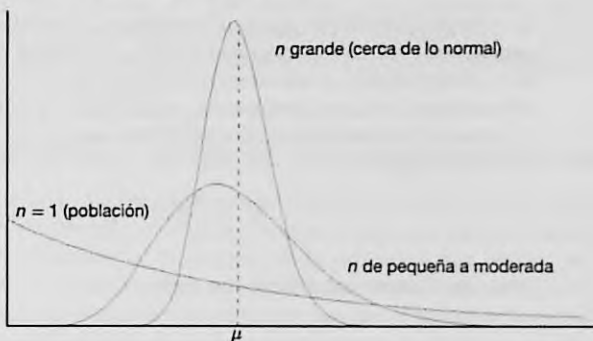


Figura 8.1: Ejemplo del teorema del límite central (distribución de \bar{X} para $n = 1$, n moderada y n grande).

Ejemplo 8.4: Una empresa de material eléctrico fabrica bombillas que tienen una duración que se distribuye aproximadamente en forma normal, con media de 800 horas y desviación estándar de 40 horas. Calcule la probabilidad de que una muestra aleatoria de 16 bombillas tenga una vida promedio de menos de 775 horas.

Solución: La distribución muestral de \bar{X} será aproximadamente normal, con $\mu_{\bar{X}} = 800$ y $\sigma_{\bar{X}} = 40/\sqrt{16} = 10$. La probabilidad que se desea es determinada por el área de la región sombreada de la figura 8.2.

En lo que corresponde a $\bar{x} = 775$, obtenemos que

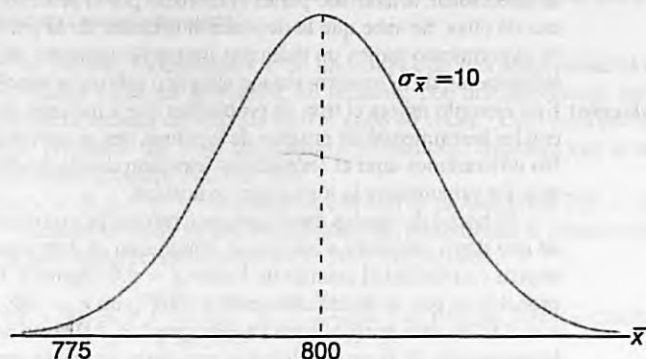


Figura 8.2: Área para el ejemplo 8.4.

$$z = \frac{775 - 800}{10} = -2.5,$$

y, por lo tanto,

$$P(\bar{X} < 775) = P(Z < -2.5) = 0.0062. \quad \blacksquare$$

Inferencias sobre la media de la población

Una aplicación muy importante del teorema del límite central consiste en determinar valores razonables de la media de la población μ . Temas como prueba de hipótesis, estimación, control de calidad y muchos otros utilizan el teorema del límite central. El siguiente ejemplo ilustra cómo se utiliza el teorema del límite central con respecto a su relación con μ , la media poblacional, aunque la aplicación formal de los temas precedentes se deja para capítulos posteriores.

En el siguiente estudio de caso proporcionamos un ejemplo en el que se hace una inferencia utilizando la distribución muestral de \bar{X} . En este ejemplo sencillo se conocen μ y σ . El teorema del límite central y el concepto general de las distribuciones muestrales a menudo se utilizan para proporcionar evidencias acerca de algún aspecto importante de una distribución, por ejemplo uno de sus parámetros. En el caso del teorema del límite central el parámetro que nos interesa es la media μ . La inferencia que se hace acerca de μ puede adoptar una de varias formas. Con frecuencia el analista desea que los datos (en la forma de \bar{x}) respalden (o no) alguna conjetura predeterminada respecto al valor de μ . El uso de lo que sabemos sobre la distribución de muestreo puede contribuir a responder este tipo de pregunta. En el siguiente estudio de caso el concepto de prueba de hipótesis conduce a un objetivo formal que destacaremos en capítulos posteriores.

Estudio de caso 8.1: Partes para automóviles. Un importante proceso de fabricación produce partes de componentes cilíndricos para la industria automotriz. Es importante que el proceso produzca partes que tengan un diámetro medio de 5.0 milímetros. El ingeniero implicado asume

que la media de la población es de 5.0 milímetros. Se lleva a cabo un experimento donde se seleccionan al azar 100 partes elaboradas por el proceso y se mide el diámetro de cada una de ellas. Se sabe que la desviación estándar de la población es $\sigma = 0.1$ milímetros. El experimento indica un diámetro promedio muestral de $\bar{x} = 5.027$ milímetros. ¿Esta información de la muestra parece apoyar o refutar la suposición del ingeniero?

Solución: Este ejemplo refleja el tipo de problemas que a menudo se presentan y que se resuelven con las herramientas de pruebas de hipótesis que se presentan en los siguientes capítulos. No utilizaremos aquí el formalismo asociado con la prueba de hipótesis, pero ilustraremos los principios y la lógica que se utilizan.

El hecho de que los datos apoyen o refuten la suposición depende de la probabilidad de que datos similares a los que se obtuvieron en este experimento ($\bar{x} = 5.027$) pueden ocurrir con facilidad cuando de hecho $\mu = 5.0$ (figura 8.3). En otras palabras, ¿qué tan probable es que se pueda obtener $\bar{x} \geq 5.027$ con $n = 100$, si la media de la población es $\mu = 5.0$? Si esta probabilidad sugiere que $\bar{x} = 5.027$ no es poco razonable, no se refuta la suposición. Si la probabilidad es muy baja, se puede argumentar con certidumbre que los datos no apoyan la suposición de que $\mu = 5.0$. La probabilidad que elegimos para el cálculo es dada por $P(|\bar{X} - 5| \geq 0.027)$.

En otras palabras, si la media μ es 5, ¿cuál es la probabilidad de que \bar{X} se desvíe

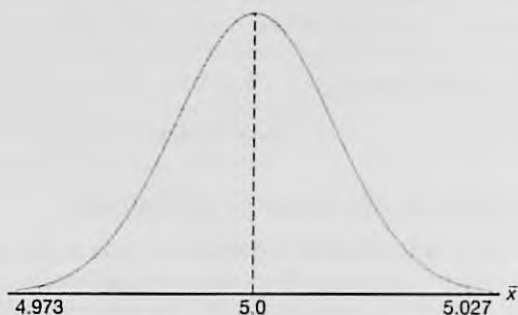


Figura 8.3: Área para el estudio de caso 8.1.

cuando mucho hasta 0.027 milímetros?

$$\begin{aligned} P(|\bar{X} - 5| \geq 0.027) &= P(\bar{X} - 5 \geq 0.027) + P(\bar{X} - 5 \leq -0.027) \\ &= 2P\left(\frac{\bar{X} - 5}{0.1/\sqrt{100}} \geq 2.7\right). \end{aligned}$$

Aquí simplemente estandarizamos \bar{X} de acuerdo con el teorema del límite central. Si la suposición $\mu = 5.0$ es cierta, $\frac{\bar{X} - 5}{0.1/\sqrt{100}}$ debería ser $N(0, 1)$. Por consiguiente,

$$2P\left(\frac{\bar{X} - 5}{0.1/\sqrt{100}} \geq 2.7\right) = 2P(Z \geq 2.7) = 2(0.0035) = 0.007.$$

Por lo tanto, se experimentaría por casualidad que una \bar{x} estaría a 0.027 milímetros

de la media en tan sólo 7 de 1000 experimentos. Como resultado, este experimento con $\bar{x} = 5.027$ ciertamente no ofrece evidencia que apoye la suposición de que $\mu = 5.0$. De hecho, ¡la refuta consistentemente! ─

Ejemplo 8.5: El viaje en un autobús especial para ir de un campus de una universidad al campus de otra en una ciudad toma, en promedio, 28 minutos, con una desviación estándar de 5 minutos. En cierta semana un autobús hizo el viaje 40 veces. ¿Cuál es la probabilidad de que el tiempo promedio del viaje sea mayor a 30 minutos? Suponga que el tiempo promedio se redondea al entero más cercano.

Solución: En este caso $\mu = 28$ y $\sigma = 5$. Necesitamos calcular la probabilidad $P(\bar{X} > 30)$ con $n = 40$. Como el tiempo se mide en una escala continua redondeada al minuto más cercano, una \bar{x} mayor que 30 sería equivalente a $\bar{x} \geq 30.5$. Por lo tanto,

$$P(\bar{X} > 30) = P\left(\frac{\bar{X} - 28}{5/\sqrt{40}} \geq \frac{30.5 - 28}{5/\sqrt{40}}\right) = P(Z \geq 3.16) = 0.0008.$$

Hay sólo una ligera probabilidad de que el tiempo promedio de un viaje del autobús exceda 30 minutos. En la figura 8.4 se presenta una gráfica ilustrativa. ─



Figura 8.4: Área para el ejemplo 8.5.

Distribución muestral de la diferencia entre dos medias

La ilustración del estudio de caso 8.1 se refiere a conceptos de inferencia estadística sobre una sola media μ . El ingeniero estaba interesado en respaldar una suposición con respecto a una sola media de población. Una aplicación mucho más importante incluye dos poblaciones. Un científico o ingeniero se podrían interesar en un experimento donde se comparan dos métodos de producción: el 1 y el 2. La base para tal comparación es $\mu_1 - \mu_2$, la diferencia entre las medias de población.

Suponga que tenemos dos poblaciones, la primera con media μ_1 y varianza σ_1^2 , y la segunda con media μ_2 y varianza σ_2^2 . Representemos con el estadístico \bar{X}_1 la media

de una muestra aleatoria de tamaño n_1 , seleccionada de la primera población, y con el estadístico \bar{X}_2 la media de una muestra aleatoria de tamaño n_2 , seleccionada de la segunda población, independiente de la muestra de la primera población. ¿Qué podríamos decir acerca de la distribución muestral de la diferencia $\bar{X}_1 - \bar{X}_2$ para muestras repetidas de tamaños n_1 y n_2 ? De acuerdo con el teorema 8.2, tanto la variable \bar{X}_1 como la variable \bar{X}_2 están distribuidas más o menos de forma normal con medias μ_1 y μ_2 y varianzas σ_1^2/n_1 y σ_2^2/n_2 , respectivamente. Esta aproximación mejora a medida que aumentan n_1 y n_2 . Al elegir muestras independientes de las dos poblaciones nos aseguramos de que las variables \bar{X}_1 y \bar{X}_2 sean independientes y, usando el teorema 7.11, con $a_1 = 1$ y $a_2 = -1$, concluimos que $\bar{X}_1 - \bar{X}_2$ se distribuye aproximadamente de forma normal con media

$$\mu_{\bar{X}_1 - \bar{X}_2} = \mu_{\bar{X}_1} - \mu_{\bar{X}_2} = \mu_1 - \mu_2$$

y varianza

$$\sigma_{\bar{X}_1 - \bar{X}_2}^2 = \sigma_{\bar{X}_1}^2 + \sigma_{\bar{X}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}.$$

El teorema del límite central se puede ampliar fácilmente al caso de dos muestras y dos poblaciones.

Teorema 8.3: Si se extraen al azar muestras independientes de tamaños n_1 y n_2 de dos poblaciones, discretas o continuas, con medias μ_1 y μ_2 y varianzas σ_1^2 y σ_2^2 , respectivamente, entonces la distribución muestral de las diferencias de las medias, $\bar{X}_1 - \bar{X}_2$, tiene una distribución aproximadamente normal, con media y varianza dadas por

$$\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2 \text{ y } \sigma_{\bar{X}_1 - \bar{X}_2}^2 = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}.$$

De aquí,

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{(\sigma_1^2/n_1) + (\sigma_2^2/n_2)}}$$

es aproximadamente una variable normal estándar.

Si tanto n_1 como n_2 son mayores o iguales que 30, la aproximación normal para la distribución de $\bar{X}_1 - \bar{X}_2$ es muy buena cuando las distribuciones subyacentes no están tan alejadas de la normal. Sin embargo, aun cuando n_1 y n_2 sean menores que 30, la aproximación normal es hasta cierto punto buena, excepto cuando las poblaciones no son definitivamente normales. Por supuesto, si ambas poblaciones son normales, entonces $\bar{X}_1 - \bar{X}_2$ tiene una distribución normal sin importar de qué tamaño sean n_1 y n_2 .

La utilidad de la distribución muestral de la diferencia entre los dos promedios muestrales es muy similar a la que se describe en el estudio de caso 8.1 en la página 235 para el caso de una sola media. Ahora presentaremos el estudio de caso 8.2, que se enfoca en el uso de la diferencia entre dos medias muestrales para respaldar (o no) la suposición de que dos medias de población son iguales.

Estudio de caso 8.2: Tiempo de secado de pinturas. Se llevan a cabo dos experimentos independientes en los que se comparan dos tipos diferentes de pintura, el A y el B. Con la pintura tipo A se pintan 18 especímenes y se registra el tiempo (en horas) que cada uno tarda en secar. Lo mismo se hace con la pintura tipo B. Se sabe que la desviación estándar de población de ambas es 1.0.

Si se supone que los especímenes pintados se secan en el mismo tiempo medio con los dos tipos de pintura, calcule $P(\bar{X}_A - \bar{X}_B > 1.0)$, donde \bar{X}_A y \bar{X}_B son los tiempos promedio de secado para muestras de tamaño $n_A = n_B = 18$.

Solución: A partir de la distribución de muestreo de $\bar{X}_A - \bar{X}_B$ sabemos que la distribución es aproximadamente normal con media

$$\mu_{\bar{X}_A - \bar{X}_B} = \mu_A - \mu_B = 0$$

y varianza

$$\sigma_{\bar{X}_A - \bar{X}_B}^2 = \frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B} = \frac{1}{18} + \frac{1}{18} = \frac{1}{9}.$$

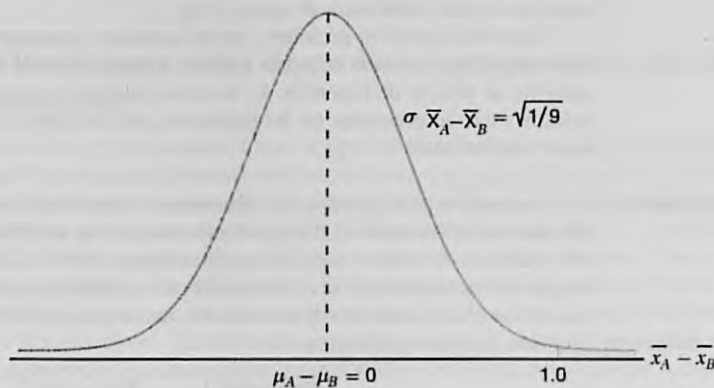


Figura 8.5: Área para el estudio de caso 8.2.

La probabilidad que se desea es dada por la región sombreada en la figura 8.5. En correspondencia con el valor $\bar{X}_A - \bar{X}_B = 1.0$, tenemos

$$z = \frac{1 - (\mu_A - \mu_B)}{\sqrt{1/9}} = \frac{1 - 0}{\sqrt{1/9}} = 3.0;$$

de modo que

$$P(Z > 3.0) = 1 - P(Z < 3.0) = 1 - 0.9987 = 0.0013. \quad \blacksquare$$

¿Qué aprendemos del estudio de caso 8.2?

La mecánica en el cálculo se basa en la suposición de que $\mu_A = \mu_B$. Suponga, sin embargo, que el experimento realmente se lleva a cabo con el fin de hacer una inferencia respecto a la igualdad de μ_A y μ_B , los tiempos medios de secado de las dos poblaciones. Si se encontrara que los dos promedios difieren por una hora (o más), este resultado sería una evidencia que nos llevaría a concluir que el tiempo medio de secado de la población

no es igual para los dos tipos de pintura. Por otro lado, suponga que la diferencia en los dos promedios muestrales es tan pequeña como, digamos, 15 minutos. Si $\mu_A = \mu_B$,

$$\begin{aligned} P[(\bar{X}_A - \bar{X}_B) > 0.25 \text{ horas}] &= P\left(\frac{\bar{X}_A - \bar{X}_B - 0}{\sqrt{1/9}} > \frac{3}{4}\right) \\ &= P\left(Z > \frac{3}{4}\right) = 1 - P(Z < 0.75) = 1 - 0.7734 = 0.2266. \end{aligned}$$

Como esta probabilidad no es baja, se concluiría que una diferencia de 15 minutos en las medias de las muestras puede ocurrir por azar, es decir, sucede con frecuencia aunque $\mu_A = \mu_B$. Por lo tanto, este tipo de diferencia en el tiempo promedio de secado ciertamente *no es una señal clara* de que $\mu_A \neq \mu_B$.

Como indicamos al principio, en los capítulos siguientes se observará un formalismo más detallado con respecto a éste y a otros tipos de inferencia estadística, por ejemplo, la prueba de hipótesis. El teorema del límite central y las distribuciones de muestreo que se presentan en las siguientes tres secciones también desempeñarán un papel fundamental.

Ejemplo 8.6: Los cinescopios para televisor del fabricante *A* tienen una duración media de 6.5 años y una desviación estándar de 0.9 años; mientras que los del fabricante *B* tienen una duración media de 6.0 años y una desviación estándar de 0.8 años. ¿Cuál es la probabilidad de que una muestra aleatoria de 36 cinescopios del fabricante *A* tenga por lo menos 1 año más de vida media que una muestra de 49 cinescopios del fabricante *B*?

Solución: Tenemos la siguiente información:

Población 1	Población 2
$\mu_1 = 6.5$	$\mu_2 = 6.0$
$\sigma_1 = 0.9$	$\sigma_2 = 0.8$
$n_1 = 36$	$n_2 = 49$

Si utilizamos el teorema 8.3, la distribución muestral de $\bar{X}_1 - \bar{X}_2$ será aproximadamente normal y tendrá una media y una desviación estándar de

$$\mu_{\bar{X}_1 - \bar{X}_2} = 6.5 - 6.0 = 0.5 \quad \text{y} \quad \sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{0.81}{36} + \frac{0.64}{49}} = 0.189.$$

La probabilidad de que 36 cinescopios del fabricante *A* tengan por lo menos 1 año más de vida media que 49 cinescopios del fabricante *B* es dada por el área de la región sombreada de la figura 8.6. Con respecto al valor $\bar{x}_1 - \bar{x}_2 = 1.0$, encontramos que

$$z = \frac{1.0 - 0.5}{0.189} = 2.65.$$

y de aquí

$$\begin{aligned} P(\bar{X}_1 - \bar{X}_2 \geq 1.0) &= P(Z > 2.65) = 1 - P(Z < 2.65) \\ &= 1 - 0.9960 = 0.0040. \end{aligned}$$

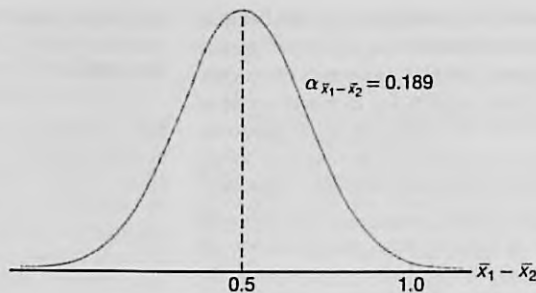


Figura 8.6: Área para el ejemplo 8.6.

Más sobre la distribución muestral de medias. Aproximación normal a la distribución binomial

En la sección 6.5 analizamos a fondo la aproximación normal a la distribución binomial. Estaban dadas las condiciones sobre los parámetros n y p , para los cuales la distribución de una variable aleatoria binomial se puede aproximar mediante la distribución normal. Los ejemplos y los ejercicios reflejaron la importancia del concepto de “aproximación normal”. Resulta que el teorema del límite central da más idea de cómo y por qué funciona esta aproximación. Sabemos con certeza que una variable aleatoria binomial es el número X de éxitos en n pruebas independientes, donde el resultado de cada prueba es binario. En el capítulo 1 también vimos que la proporción calculada en un experimento así es un promedio de un conjunto de ceros y unos. De hecho, mientras que la proporción X/n es un promedio, X es la suma de este conjunto de ceros y unos, y tanto X como X/n son casi normales si n es suficientemente grande. Desde luego, a partir de lo que aprendimos en el capítulo 6, sabemos que hay condiciones de n y p que afectan la calidad de la aproximación; a saber, $np \geq 5$ y $nq \geq 5$.

Ejercicios

8.17 Si se extraen todas las muestras posibles de tamaño 16 de una población normal con media igual a 50 y desviación estándar igual a 5, ¿cuál es la probabilidad de que una media muestral \bar{X} caiga en el intervalo que va de $\mu_{\bar{X}} - 1.9\sigma_{\bar{X}}$ a $\mu_{\bar{X}} - 0.4\sigma_{\bar{X}}$? Suponga que las medias muestrales se pueden medir con cualquier grado de precisión.

8.18 Si la desviación estándar de la media para la distribución muestral de muestras aleatorias de tamaño 36 de una población grande o infinita es 2, ¿qué tan grande debe ser el tamaño de la muestra si la desviación estándar se reduce a 1.2?

8.19 Se fabrica cierto tipo de hilo con una resistencia a la tensión media de 78.3 kilogramos y una desviación estándar de 5.6 kilogramos. ¿Cómo cambia la varianza de la media muestral cuando el tamaño de la muestra

- a) aumenta de 64 a 196?
b) disminuye de 784 a 49?

8.20 Dada la población uniforme discreta

$$f(x) = \begin{cases} \frac{1}{3}, & x = 2, 4, 6, \\ 0, & \text{en otro caso,} \end{cases}$$

calcule la probabilidad de que una muestra aleatoria de tamaño 54, seleccionada con reemplazo, produzca una media muestral mayor que 4.1 pero menor que 4.4. Suponga que las medias se miden al décimo más cercano.

8.21 Una máquina de bebidas gaseosas se ajusta de manera que la cantidad de bebida que sirve promedie 240 mililitros con una desviación estándar de 15 mililitros. La máquina se verifica periódicamente tomando una muestra de 40 bebidas y calculando el

contenido promedio. Si la media de las 40 bebidas es un valor dentro del intervalo $\mu_{\bar{x}} \pm 2\sigma_{\bar{x}}$, se piensa que la máquina opera satisfactoriamente; de lo contrario, se ajusta. En la sección 8.3 el ejecutivo de la empresa encontró que la media de 40 bebidas era $\bar{x} = 236$ mililitros y concluyó que la máquina no necesitaba un ajuste. ¿Fue ésta una decisión razonable?

8.22 Las estaturas de 1000 estudiantes se distribuyen aproximadamente de forma normal con una media de 174.5 centímetros y una desviación estándar de 6.9 centímetros. Si se extraen 200 muestras aleatorias de tamaño 25 de esta población y las medias se registran al décimo de centímetro más cercano, determine

- la media y la desviación estándar de la distribución muestral de \bar{X} ;
- el número de las medias muestrales que caen entre 172.5 y 175.8 centímetros;
- el número de medias muestrales que caen por debajo de 172.0 centímetros.

8.23 La variable aleatoria X , que representa el número de cerezas en un tarta, tiene la siguiente distribución de probabilidad:

x	4	5	6	7
$P(X = x)$	0.2	0.4	0.3	0.1

- Calcule la media μ y la varianza σ^2 de X .
- Calcule la media $\mu_{\bar{X}}$ y la varianza $\sigma_{\bar{X}}^2$ de la media \bar{X} para muestras aleatorias de 36 tartas de cereza.
- Calcule la probabilidad de que el número promedio de cerezas en 36 tartas sea menor que 5.5.

8.24 Si cierta máquina fabrica resistencias eléctricas que tienen una resistencia media de 40 ohms y una desviación estándar de 2 ohms, ¿cuál es la probabilidad de que una muestra aleatoria de 36 de estas resistencias tenga una resistencia combinada de más de 1458 ohms?

8.25 La vida media de una máquina para elaborar panes de 7 años, con una desviación estándar de 1 año. Suponga que la vida de estas máquinas sigue aproximadamente una distribución normal y calcule

- la probabilidad de que la vida media de una muestra aleatoria de 9 de estas máquinas caiga entre 6.4 y 7.2 años;
- el valor de x a la derecha del cual caería 15% de las medias calculadas de muestras aleatorias de tamaño 9.

8.26 La cantidad de tiempo que le toma al cajero de un banco con servicio en el automóvil atender a un cliente es una variable aleatoria con una media $\mu = 3.2$ minutos y una desviación estándar $\sigma = 1.6$ minutos. Si se observa una muestra aleatoria de 64 clientes, calcule la probabilidad de que el tiempo medio que el cliente

pasa en la ventanilla del cajero sea

- a lo sumo 2.7 minutos;
- más de 3.5 minutos;
- al menos 3.2 minutos pero menos de 3.4 minutos.

8.27 En un proceso químico la cantidad de cierto tipo de impureza en el producto es difícil de controlar y por ello es una variable aleatoria. Se especula que la cantidad media de la población de impurezas es 0.20 gramos por gramo del producto. Se sabe que la desviación estándar es 0.1 gramos por gramo. Se realiza un experimento para entender mejor la especulación de que $\mu = 0.2$. El proceso se lleva a cabo 50 veces en un laboratorio y el promedio de la muestra \bar{x} resulta ser 0.23 gramos por gramo. Comente sobre la especulación de que la cantidad media de impurezas es 0.20 gramos por gramo. Utilice el teorema del límite central en su respuesta.

8.28 Se toma una muestra aleatoria de tamaño 25 de una población normal que tiene una media de 80 y una desviación estándar de 5. Una segunda muestra aleatoria de tamaño 36 se toma de una población normal diferente que tiene una media de 75 y una desviación estándar de 3. Calcule la probabilidad de que la media muestral calculada de las 25 mediciones exceda la media muestral calculada de las 36 mediciones por lo menos 3.4 pero menos de 5.9. Suponga que las diferencias de las medias se miden al décimo más cercano.

8.29 La distribución de alturas de cierta raza de perros *terrier* tiene una media de 72 centímetros y una desviación estándar de 10 centímetros; en tanto que la distribución de alturas de cierta raza de *poodles* tiene una media de 28 centímetros con una desviación estándar de 5 centímetros. Suponga que las medias muestrales se pueden medir con cualquier grado de precisión y calcule la probabilidad de que la media muestral de una muestra aleatoria de alturas de 64 *terriers* exceda la media muestral para una muestra aleatoria de alturas de 100 *poodles* a lo sumo 44.2 centímetros.

8.30 La calificación promedio de los estudiantes de primer año en un examen de aptitudes en cierta universidad es 540, con una desviación estándar de 50. Suponga que las medias se miden con cualquier grado de precisión. ¿Cuál es la probabilidad de que dos grupos seleccionados al azar, que constan de 32 y 50 estudiantes, respectivamente, difieran en sus calificaciones promedio por

- más de 20 puntos?
- una cantidad entre 5 y 10 puntos?

8.31 Considere el estudio de caso 8.2 de la página 238. Suponga que en un experimento se utilizaron 18 especímenes para cada tipo de pintura y que $\bar{x}_A - \bar{x}_B$ la diferencia real en el tiempo medio de secado, resultó ser 1.0.

- a) ¿Parecería ser un resultado razonable si los dos tiempos promedio de secado de las dos poblaciones realmente son iguales? Utilice el resultado que se obtuvo en el estudio de caso 8.2.
- b) Si alguien hiciera el experimento 10,000 veces bajo la condición de que $\mu_A = \mu_B$, ¿en cuántos de esos 10,000 experimentos habría una diferencia $\bar{x}_A - \bar{x}_B$ tan grande como 1.0 (o más grande)?

8.32 Dos máquinas diferentes de llenado de cajas se utilizan para llenar cajas de cereal en una línea de ensamble. La medición fundamental en la que influyen estas máquinas es el peso del producto en las cajas. Los ingenieros están seguros de que la varianza en el peso del producto es $\sigma^2 = 1$ onza. Se realizan experimentos usando ambas máquinas con tamaños muestrales de 36 cada una. Los promedios muestrales para las máquinas A y B son $\bar{x}_A = 4.5$ onzas y $\bar{x}_B = 4.7$ onzas. Los ingenieros se sorprenden de que los dos promedios muestrales para las máquinas de llenado sean tan diferentes.

- a) Utilice el teorema del límite central para determinar

$$P(\bar{X}_B - \bar{X}_A \geq 0.2)$$

bajo la condición de que $\mu_A = \mu_B$.

- b) ¿Los experimentos mencionados parecen, de cualquier forma, apoyar consistentemente la suposición de que las medias de población de las dos máquinas son diferentes? Explique utilizando la respuesta que encontró en el inciso a.

8.33 El benceno es una sustancia química altamente tóxica para los seres humanos. Sin embargo, se utiliza en la fabricación de medicamentos, de tintes y de recubrimientos, así como en la peletería. Las regulaciones del gobierno establecen que el contenido de benceno en el agua que resulte de cualquier proceso de producción en el que participe esta sustancia no debe exceder 7950 partes por millón (ppm). Para un proceso particular de interés, un fabricante recolectó una muestra de agua 25 veces de manera aleatoria y el promedio muestral \bar{x} fue de 7960 ppm. A partir de los datos históricos, se sabe que la desviación estándar σ es 100 ppm.

- a) ¿Cuál es la probabilidad de que el promedio muestral en este experimento exceda el límite establecido por el gobierno, si la media de la población es igual al límite? Utilice el teorema del límite central.
- b) ¿La $\bar{x} = 7960$ observada en este experimento es firme evidencia de que la media de la población

en este proceso excede el límite impuesto por el gobierno? Responda calculando

$$P(\bar{X} \geq 7960 \mid \mu = 7950).$$

Suponga que la distribución de la concentración de benceno es normal.

8.34 En la fabricación de cierto producto de acero se están utilizando dos aleaciones, la A y la B . Se necesita diseñar un experimento para comparar las dos aleaciones en términos de su capacidad de carga máxima en toneladas, es decir, la cantidad máxima de carga que pueden soportar sin romperse. Se sabe que las dos desviaciones estándar de la capacidad de carga son iguales a 5 toneladas cada una. Se realiza un experimento en el que se prueban 30 especímenes de cada aleación (A y B) y se obtienen los siguientes resultados:

$$\bar{x}_A = 49.5, \quad \bar{x}_B = 45.5; \quad \bar{x}_A - \bar{x}_B = 4.$$

Los fabricantes de la aleación A están convencidos de que esta evidencia demuestra de forma concluyente que $\mu_A > \mu_B$ y, por lo tanto, que su aleación es mejor. Los fabricantes de la aleación B afirman que el experimento fácilmente podría haber resultado $\bar{x}_A - \bar{x}_B = 4$, *incluso si* las dos medias de población fueran iguales. En otras palabras, "¡los resultados no son concluyentes!".

- a) Encuentre un argumento que ponga en evidencia el error de los fabricantes de la aleación B . Para ello calcule

$$P(\bar{X}_A - \bar{X}_B > 4 \mid \mu_A = \mu_B).$$

- b) ¿Considera que estos datos apoyan fuertemente a la aleación A ?

8.35 Considere la situación del ejemplo 8.4 de la página 234. ¿Los resultados que se obtuvieron allí lo llevan a cuestionar la premisa de que $\mu = 800$ horas? Proporcione un resultado probabilístico que indique qué tan raro es el evento $\bar{X} \leq 775$ cuando $\mu = 800$. Por otro lado, ¿qué tan raro sería si μ fuera, verdaderamente, digamos, $\neq 760$ horas?

8.36 Sea X_1, X_2, \dots, X_n una muestra aleatoria de una distribución que sólo puede adoptar valores positivos. Utilice el teorema del límite central para argumentar que si n es tan grande como se requiere, entonces $Y = X_1 X_2 \dots X_n$ tiene aproximadamente una distribución logarítmica normal.

8.5 Distribución muestral de S^2

En la sección anterior aprendimos acerca de la distribución muestral de \bar{X} . El teorema del límite central nos permitió utilizar el hecho de que

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

tiende a $N(0, 1)$ a medida que crece el tamaño de la muestra. *Las distribuciones muestrales de estadísticos importantes* nos permiten conocer información sobre los parámetros. Por lo general, los parámetros son las contrapartes del estadístico en cuestión. Por ejemplo, si un ingeniero se interesa en la resistencia media de la población de cierto tipo de resistencia, sacará provecho de la distribución muestral de \bar{X} una vez que reúna la información de la muestra. Por otro lado, si está estudiando la variabilidad en la resistencia, evidentemente utilizará la distribución muestral de S^2 para conocer la contraparte paramétrica, la varianza de la población σ^2 .

Si se extrae una muestra aleatoria de tamaño n de una población normal con media μ y varianza σ^2 , y se calcula la varianza muestral, se obtiene un valor del estadístico S^2 . Procederemos a considerar la distribución del estadístico $(n-1)S^2/\sigma^2$.

Mediante la suma y la resta de la media muestral \bar{X} es fácil ver que

$$\begin{aligned}\sum_{i=1}^n (X_i - \mu)^2 &= \sum_{i=1}^n [(X_i - \bar{X}) + (\bar{X} - \mu)]^2 \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n (\bar{X} - \mu)^2 + 2(\bar{X} - \mu) \sum_{i=1}^n (X_i - \bar{X}) \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + n(\bar{X} - \mu)^2.\end{aligned}$$

Al dividir cada término de la igualdad entre σ^2 y sustituir $(n-1)S^2$ por $\sum_{i=1}^n (X_i - \bar{X})^2$, obtenemos

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 = \frac{(n-1)S^2}{\sigma^2} + \frac{(\bar{X} - \mu)^2}{\sigma^2/n}.$$

Ahora, de acuerdo con el corolario 7.1 de la página 222, sabemos que

$$\sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^2}$$

es una variable aleatoria chi cuadrada con n grados de libertad. Tenemos una variable aleatoria chi cuadrada con n grados de libertad dividida en dos componentes. Observe que en la sección 6.7 demostramos que una distribución chi cuadrada es un caso especial de la distribución gamma. El segundo término del lado derecho es Z^2 , que es una variable aleatoria chi cuadrada con 1 grado de libertad, y resulta que $(n-1)S^2/\sigma^2$ es una variable aleatoria chi cuadrada con $n-1$ grados de libertad. Formalizamos esto en el siguiente teorema.

Teorema 8.4: Si S^2 es la varianza de una muestra aleatoria de tamaño n que se toma de una población normal que tiene la varianza σ^2 , entonces el estadístico

$$\chi^2 = \frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2}$$

tiene una distribución chi cuadrada con $v = n-1$ grados de libertad.

Los valores de la variable aleatoria χ^2 se calculan de cada muestra mediante la fórmula

$$\chi^2 = \frac{(n-1)s^2}{\sigma^2}.$$

La probabilidad de que una muestra aleatoria produzca un valor χ^2 mayor que algún valor específico es igual al área bajo la curva a la derecha de este valor. El valor χ^2 por arriba del cual se encuentra un área de α por lo general se representa con χ^2_α . Esto se ilustra mediante la región sombreada de la figura 8.7.

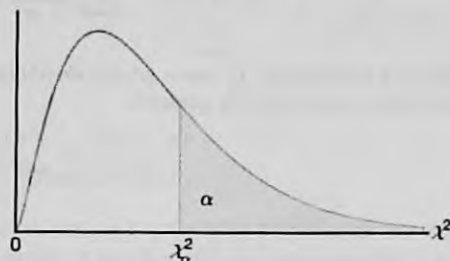


Figura 8.7: La distribución chi cuadrada.

La tabla A.5 da los valores de χ^2_α para diversos valores de α y ν . Las áreas, α , son los encabezados de las columnas; los grados de libertad, ν , se dan en la columna izquierda, y las entradas de la tabla son los valores χ^2 . En consecuencia, el valor χ^2 con 7 grados de libertad, que deja un área de 0.05 a la derecha, es $\chi^2_{0.05} = 14.067$. Debido a la falta de simetría, para encontrar $\chi^2_{0.95} = 2.167$ para $\nu = 7$ también debemos usar las tablas.

Exactamente 95% de una distribución chi cuadrada cae entre $\chi^2_{0.975}$ y $\chi^2_{0.025}$. Un valor χ^2 que cae a la derecha de $\chi^2_{0.025}$ no tiene probabilidades de ocurrir, a menos que el valor de σ^2 que supusimos sea demasiado pequeño. Lo mismo sucede con un valor χ^2 que cae a la izquierda de $\chi^2_{0.975}$, el cual tampoco es probable que ocurra, a menos que el valor de σ^2 que supusimos sea demasiado grande. En otras palabras, es posible tener un valor χ^2 a la izquierda de $\chi^2_{0.975}$ o a la derecha de $\chi^2_{0.025}$ cuando el valor de σ^2 es correcto; pero si esto sucediera, lo más probable es que el valor de σ^2 que se supuso sea un error.

Ejemplo 8.7: Un fabricante de baterías para automóvil garantiza que su producto durará, en promedio, 3 años con una desviación estándar de 1 año. Si cinco de estas baterías tienen duraciones de 1.9, 2.4, 3.0, 3.5 y 4.2 años, ¿el fabricante continuará convencido de que sus baterías tienen una desviación estándar de 1 año? Suponga que las duraciones de las baterías siguen una distribución normal.

Solución: Primero se calcula la varianza de la muestra usando el teorema 8.1,

$$s^2 = \frac{(5)(48.26) - (15)^2}{(5)(4)} = 0.815.$$

Entonces,

$$\chi^2 = \frac{(4)(0.815)}{1} = 3.26$$

es un valor de una distribución chi cuadrada con 4 grados de libertad. Como 95% de los valores χ^2 con 4 grados de libertad cae entre 0.484 y 11.143, el valor calculado con $\sigma^2 = 1$ es razonable y, por lo tanto, el fabricante no tiene razones para sospechar que la desviación estándar no sea igual a 1 año.

Grados de libertad como una medición de la información muestral

Del corolario 7.1 expuesto en la sección 7.3 recuerde que

$$\sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^2}$$

tiene una distribución χ^2 con n grados de libertad. Observe también el teorema 8.4, el cual indica que la variable aleatoria

$$\frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2}$$

tiene una distribución χ^2 con $n-1$ grados de libertad. El lector debe también recordar que el término *grados de libertad*, que se utiliza en este contexto idéntico, se estudió en el capítulo 1.

Como antes indicamos, el teorema 8.4 no se demostrará; sin embargo, el lector puede verlo como una indicación de que cuando no se conoce μ y se considera la distribución de

$$\sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2},$$

hay **1 grado menos de libertad**, o se pierde un grado de libertad al estimar μ (es decir, cuando μ se reemplaza por \bar{x}). En otras palabras, en la muestra aleatoria de la distribución normal hay n grados de libertad o *partes de información* independientes. Cuando los datos (los valores en la muestra) se utilizan para calcular la media, hay un grado menos de libertad en la información que se utiliza para estimar σ^2 .

8.6 Distribución t

En la sección 8.4 se analizó la utilidad del teorema del límite central. Sus aplicaciones giran en torno a las inferencias sobre una media de la población o a la diferencia entre dos medias de población. En este contexto es evidente la utilidad de utilizar el teorema del límite central y la distribución normal. Sin embargo, se supuso que se conoce la desviación estándar de la población. Esta suposición quizá sea razonable en situaciones en las que el ingeniero está muy familiarizado con el sistema o proceso. Sin embargo, en muchos escenarios experimentales el conocimiento de σ no es ciertamente más razonable que el conocimiento de la media de la población μ . A menudo, de hecho, una estimación de σ debe ser proporcionada por la misma información muestral que produce el promedio muestral \bar{x} . Como resultado, un estadístico natural a considerar para tratar con las inferencias sobre μ es

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}},$$

dado que S es el análogo de la muestra para σ . Si el tamaño de la muestra es pequeño, los valores de S^2 fluctúan de forma considerable de una muestra a otra (véase el ejercicio 8.43 de la página 259) y la distribución de T se desvía de forma apreciable de la de una distribución normal estándar.

Si el tamaño de la muestra es suficientemente grande, digamos $n \geq 30$, la distribución de T no difiere mucho de la normal estándar. Sin embargo, para $n < 30$ es útil tratar con la distribución exacta de T . Para desarrollar la distribución muestral de T , supondremos que nuestra muestra aleatoria se seleccionó de una población normal. Podemos escribir, entonces,

$$T = \frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{\sqrt{S^2/\sigma^2}} = \frac{Z}{\sqrt{V/(n-1)}},$$

donde

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

tiene una distribución normal estándar y

$$V = \frac{(n-1)S^2}{\sigma^2}$$

tiene una distribución chi cuadrada con $v = n - 1$ grados de libertad. Al obtener muestras de poblaciones normales se puede demostrar que \bar{X} y S^2 son independientes y, en consecuencia, también lo son Z y V . El siguiente teorema proporciona la definición de una variable aleatoria T como una función de Z (normal estándar) y χ^2 . Para completar se proporciona la función de densidad de la distribución t .

Teorema 8.5: Sea Z una variable aleatoria normal estándar y V una variable aleatoria chi cuadrada con v grados de libertad. Si Z y V son independientes, entonces la distribución de la variable aleatoria T , donde

$$T = \frac{Z}{\sqrt{V/v}},$$

es dada por la función de densidad

$$h(t) = \frac{\Gamma[(v+1)/2]}{\Gamma(v/2)\sqrt{\pi v}} \left(1 + \frac{t^2}{v}\right)^{-(v+1)/2}, \quad -\infty < t < \infty.$$

Ésta se conoce como la **distribución t** con v grados de libertad.

A partir de lo antes expuesto, y del teorema anterior, se deriva el siguiente corolario.

Corolario 8.1: Sean X_1, X_2, \dots, X_n variables aleatorias independientes normales con media μ y desviación estándar σ . Sea

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{y} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Entonces la variable aleatoria $\frac{\bar{X} - \mu}{s/\sqrt{n}}$ tiene una distribución t con $\nu = n - 1$ grados de libertad.

La distribución de probabilidad de T se publicó por primera vez en 1908 en un artículo de W. S. Gosset. En esa época, Gosset trabajaba para una cervecería irlandesa que prohibía a sus empleados que publicaran los resultados de sus investigaciones. Para evadir la prohibición Gosset publicó su trabajo en secreto bajo el seudónimo de "Student". Es por esto que a la distribución de T se le suele llamar distribución t de Student o simplemente distribución t . Para derivar la ecuación de esta distribución Gosset supuso que las muestras se seleccionaban de una población normal. Aunque ésta parecería una suposición muy restrictiva, se puede demostrar que las poblaciones que no son normales y que poseen distribuciones en forma casi de campana aún proporcionan valores de T que se aproximan muy de cerca a la distribución t .

¿Qué apariencia tiene la distribución t ?

La distribución de T se parece a la distribución de Z en que ambas son simétricas alrededor de una media de cero. Ambas distribuciones tienen forma de campana, pero la distribución t es más variable debido al hecho de que los valores T dependen de las fluctuaciones de dos cantidades, \bar{X} y S^2 ; mientras que los valores Z dependen sólo de los cambios en \bar{X} de una muestra a otra. La distribución de T difiere de la de Z en que la varianza de T depende del tamaño de la muestra n y siempre es mayor que 1. Sólo cuando el tamaño de la muestra $n \rightarrow \infty$ las dos distribuciones serán iguales. En la figura 8.8 se presenta la relación entre una distribución normal estándar ($\nu = \infty$) y las distribuciones t con 2 y 5 grados de libertad. Los puntos porcentuales de la distribución t se dan en la tabla A.4.



Figura 8.8: Curvas de la distribución t para $\nu = 2, 5$ y ∞ .

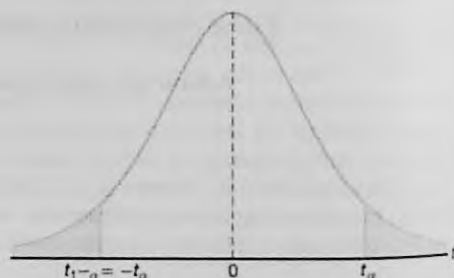


Figura 8.9: Propiedad de simetría (alrededor de 0) de la distribución t .

El valor t por arriba del cual se encuentra un área igual a α por lo general se representa con t_{α} . Por consiguiente, el valor t con 10 grados de libertad que deja una área de 0.025 a la derecha es $t = 2.228$. Como la distribución t es simétrica alrededor de una media de cero, tenemos $t_{1-\alpha} = -t_{\alpha}$; es decir, el valor t que deja una área de $1 - \alpha$ a la derecha y, por lo tanto, una área de α a la izquierda es igual al valor t negativo que deja una área de α en la cola derecha de la distribución (véase la figura 8.9). Esto es, $t_{0.95} = -t_{0.05}$, $t_{0.99} = -t_{0.01}$, etcétera.

Ejemplo 8.8: El valor t con $\nu = 14$ grados de libertad que deja una área de 0.025 a la izquierda y, por lo tanto, una área de 0.975 a la derecha, es

$$t_{0.975} = -t_{0.025} = -2.145. \quad \blacksquare$$

Ejemplo 8.9: Calcule $P(-t_{0.025} < T < t_{0.05})$.

Solución: Como $t_{0.05}$ deja una área de 0.05 a la derecha y $-t_{0.025}$ deja una área de 0.025 a la izquierda, obtenemos una área total de

$$1 - 0.05 - 0.025 = 0.925$$

entre $-t_{0.025}$ y $t_{0.05}$. En consecuencia,

$$P(-t_{0.025} < T < t_{0.05}) = 0.925. \quad \blacksquare$$

Ejemplo 8.10: Calcule k tal que $P(k < T < -1.761) = 0.045$ para una muestra aleatoria de tamaño 15 que se selecciona de una distribución normal y $\frac{\bar{X} - \mu}{s/\sqrt{n}}$.



Figura 8.10: Valores t para el ejemplo 8.10.

Solución: A partir de la tabla A.4 advertimos que 1.761 corresponde a $t_{0.05}$ cuando $\nu = 14$. Por lo tanto, $-t_{0.05} = -1.761$. Puesto que en el enunciado de probabilidad original k está a la izquierda de $-t_{0.05} = -1.761$, tenemos que $k = -t_{\alpha}$. Entonces, a partir de la figura 8.10, tenemos

$$0.045 = 0.05 - \alpha, \text{ o } \alpha = 0.005.$$

Así, de la tabla A.4 con $\nu = 14$,

$$k = -t_{0.005} = -2.977 \text{ y } P(-2.977 < T < -1.761) = 0.045. \quad \blacksquare$$

Exactamente 95% de los valores de una distribución t con $\nu = n - 1$ grados de libertad caen entre $-t_{0.025}$ y $t_{0.025}$. Por supuesto, hay otros valores t que contienen 95% de la distribución, como $-t_{0.02}$ y $t_{0.03}$, pero estos valores no aparecen en la tabla A.4 y, además, el intervalo más corto posible se obtiene eligiendo valores t que dejen exactamente la misma área en las dos colas de nuestra distribución. Un valor t que caiga por debajo de $-t_{0.025}$ o por arriba de $t_{0.025}$ tendería a hacernos creer que ha ocurrido un evento muy raro, o que quizá nuestra suposición acerca de μ es un error. Si esto ocurriera, tendríamos que tomar la decisión de que el valor de μ que supusimos es erróneo. De hecho, un valor t que cae por debajo de $-t_{0.01}$ o por arriba de $t_{0.01}$ proporcionaría incluso evidencia más sólida de que el valor de μ que supusimos es muy improbable. En el capítulo 10 se tratarán procedimientos generales para probar aseveraciones respecto al valor del parámetro μ . El siguiente ejemplo ilustra una vista preliminar del fundamento de tales procedimientos.

Ejemplo 8.11: Un ingeniero químico afirma que el rendimiento medio de la población de un cierto proceso de lotes es 500 gramos por mililitro de materia prima. Para verificar dicha afirmación muestrea 25 lotes cada mes. Si el valor t calculado cae entre $-t_{0.05}$ y $t_{0.05}$, queda satisfecho con su afirmación. ¿Qué conclusión debería sacar de una muestra que tiene una media $\bar{x} = 518$ gramos por mililitro y una desviación estándar muestral $s = 40$ gramos? Suponga que la distribución de rendimientos es aproximadamente normal.

Solución: En la tabla A.4 encontramos que $t_{0.05} = 1.711$ para 24 grados de libertad. Por lo tanto, el ingeniero quedará satisfecho con esta afirmación si una muestra de 25 lotes rinde un valor t entre -1.711 y 1.711 . Si $\mu = 500$, entonces,

$$t = \frac{518 - 500}{40/\sqrt{25}} = 2.25,$$

un valor muy superior a 1.711. La probabilidad de obtener un valor t , con $\nu = 24$, igual o mayor que 2.25, es aproximadamente 0.02. Si $\mu > 500$, el valor de t calculado de la muestra sería más razonable. Por lo tanto, es probable que el ingeniero concluya que el proceso produce un mejor producto del que pensaba. J

¿Para qué se utiliza la distribución t ?

La distribución t se usa ampliamente en problemas relacionados con inferencias acerca de la media de la población (como se ilustra en el ejemplo 8.11) o en problemas que implican muestras comparativas (es decir, en casos donde se trata de determinar si las medias de dos muestras son muy diferentes). El uso de la distribución se ampliará en los capítulos 9, 10, 11 y 12. El lector debería notar que el uso de la distribución t para el estadístico

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

requiere que X_1, X_2, \dots, X_n sean normales. El uso de la distribución t y la consideración del tamaño de la muestra no se relacionan con el teorema del límite central. El uso de la distribución normal estándar en vez de T para $n \geq 30$ sólo implica, en este caso, que S es un estimador suficientemente bueno de σ . En los siguientes capítulos la distribución t se usa con amplitud.

8.7 Distribución F

Recomendamos la distribución t en parte por su aplicación a problemas en los que hay muestreo comparativo, es decir, a problemas en que se tienen que comparar dos medias muestrales. Por ejemplo, algunos de los ejemplos que daremos en los siguientes capítulos adoptarán un método aún más formal; un ingeniero químico reúne datos de dos catalizadores, un biólogo recoge datos sobre dos medios de crecimiento o un químico reúne datos sobre dos métodos de recubrimiento de material para prevenir la corrosión. Si bien es importante que la información muestral aclare lo relacionado con dos medias de población, a menudo éste es el caso en el que comparar la variabilidad es igual de importante, si no es que más. La distribución F tiene una amplia aplicación en la comparación de varianzas muestrales y también es aplicable en problemas que implican dos o más muestras.

El estadístico F se define como el cociente de dos variables aleatorias chi cuadrada independientes, dividida cada una entre su número de grados de libertad. En consecuencia, podemos escribir

$$F = \frac{U/v_1}{V/v_2},$$

donde U y V son variables aleatorias independientes que tienen distribuciones chi cuadrada con v_1 y v_2 grados de libertad, respectivamente. Estableceremos ahora la distribución muestral de F .

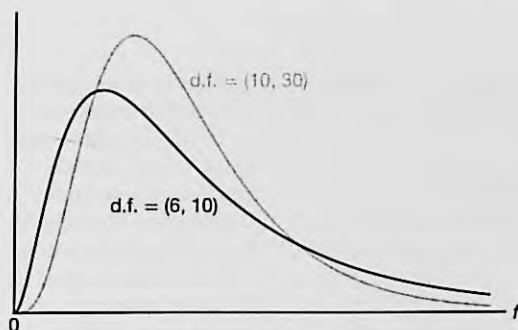
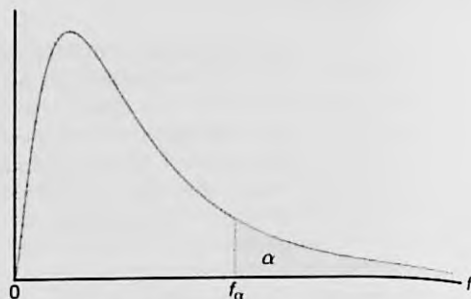
Teorema 8.6: Sean U y V dos variables aleatorias independientes que tienen distribuciones chi cuadrada con v_1 y v_2 grados de libertad, respectivamente. Entonces, la distribución de la variable aleatoria $F = \frac{U/v_1}{V/v_2}$ es dada por la función de densidad

$$h(f) = \begin{cases} \frac{\Gamma((v_1+v_2)/2)(v_1/v_2)^{v_1/2}}{\Gamma(v_1/2)\Gamma(v_2/2)} \frac{f^{(v_1/2)-1}}{(1+v_1f/v_2)^{(v_1+v_2)/2}}, & f > 0, \\ 0, & f \leq 0. \end{cases}$$

Ésta se conoce como la **distribución F** con v_1 y v_2 grados de libertad (g.l.).

En capítulos posteriores utilizaremos ampliamente la variable aleatoria F . Sin embargo, no emplearemos la función de densidad, la cual sólo se dará como complemento. La curva de la distribución F no sólo depende de los dos parámetros v_1 y v_2 sino también del orden en el que se establecen. Una vez que tenemos estos dos valores, podemos identificar la curva. En la figura 8.11 se presentan distribuciones F típicas.

Sea f_α el valor f por arriba del cual encontramos un área igual a α . Esto se ilustra mediante la región sombreada de la figura 8.12. La tabla A.6 proporciona valores de f_α sólo para $\alpha = 0.05$ y $\alpha = 0.01$ para varias combinaciones de los grados de libertad v_1 y v_2 . Por lo tanto, el valor f con 6 y 10 grados de libertad, que deja un área de 0.05 a la derecha, es $f_{0.05} = 3.22$. Por medio del siguiente teorema, la tabla A.6 también se puede utilizar para encontrar valores de $f_{0.95}$ y $f_{0.99}$. La demostración se deja al lector.

Figura 8.11: Distribuciones F típicas.Figura 8.12: Ilustración de la f_α para la distribución F .

Teorema 8.7: Al escribir $f_\alpha(v_1, v_2)$ para f_α con v_1 y v_2 grados de libertad, obtenemos

$$f_{1-\alpha}(v_1, v_2) = \frac{1}{f_\alpha(v_2, v_1)}.$$

Por consiguiente, el valor f con 6 y 10 grados de libertad, que deja una área de 0.95 a la derecha, es

$$f_{0.95}(6, 10) = \frac{1}{f_{0.05}(10, 6)} = \frac{1}{4.06} = 0.246.$$

La distribución F con dos varianzas muestrales

Suponga que las muestras aleatorias de tamaños n_1 y n_2 se seleccionan de dos poblaciones normales con varianzas σ_1^2 y σ_2^2 , respectivamente. Del teorema 8.4, sabemos que

$$\chi_1^2 = \frac{(n_1 - 1)S_1^2}{\sigma_1^2} \text{ y } \chi_2^2 = \frac{(n_2 - 1)S_2^2}{\sigma_2^2}$$

son variables aleatorias que tienen distribuciones chi cuadrada con $v_1 = n_1 - 1$ y $v_2 = n_2 - 1$ grados de libertad. Además, como las muestras se seleccionan al azar, tratamos con variables aleatorias independientes. Entonces, usando el teorema 8.6 con $\chi_1^2 = U$ y $\chi_2^2 = V$, obtenemos el siguiente resultado.

Teorema 8.8: Si S_1^2 y S_2^2 son las varianzas de muestras aleatorias independientes de tamaño n_1 y n_2 tomadas de poblaciones normales con varianzas σ_1^2 y σ_2^2 , respectivamente, entonces,

$$F = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} = \frac{\sigma_2^2 S_1^2}{\sigma_1^2 S_2^2}$$

tiene una distribución F con $v_1 = n_1 - 1$ y $v_2 = n_2 - 1$ grados de libertad.

¿Para qué se utiliza la distribución F ?

Al inicio de esta sección contestamos esta pregunta parcialmente. La distribución F se usa en situaciones de dos muestras para hacer inferencias acerca de las varianzas de población, lo cual implica aplicar el teorema 8.8. Sin embargo, la distribución F también se puede aplicar a muchos otros tipos de problemas que involucren varianzas muestrales. De hecho, la distribución F se llama *distribución de razón de varianzas*. Como ejemplo, considere el estudio de caso 8.2 en el que se compararon las dos pinturas, A y B , en relación con el tiempo medio que tardan en secar, en donde la distribución normal se aplica muy bien (suponiendo que se conocen σ_A y σ_B). Sin embargo, suponga que necesitamos comparar tres tipos de pinturas, digamos A , B y C , y que queremos determinar si las medias de población son equivalentes. Suponga que un resumen de la información importante del experimento es el siguiente:

Pintura	Media muestral	Varianza muestral	Tamaño muestral
A	$\bar{X}_A = 4.5$	$s_A^2 = 0.20$	10
B	$\bar{X}_B = 5.5$	$s_B^2 = 0.14$	10
C	$\bar{X}_C = 6.5$	$s_C^2 = 0.11$	10

El problema se centra alrededor de si los promedios muestrales (\bar{x}_A , \bar{x}_B , \bar{x}_C) están o no suficientemente alejados. La implicación de "suficientemente alejados" resulta muy importante. Parecería razonable que si la variabilidad entre los promedios muestrales es mayor que lo que se esperaría por casualidad, los datos no apoyan la conclusión de que $\mu_A = \mu_B = \mu_C$. Si estos promedios muestrales pudieran ocurrir por casualidad depende de la *variabilidad dentro de las muestras*, cuando se cuantifican por medio de s_A^2 , s_B^2 y s_C^2 . La idea de los componentes importantes de la variabilidad se observa mejor utilizando algunas gráficas sencillas. Considere la gráfica de los datos brutos de las muestras A , B y C que se presenta en la figura 8.13. Estos datos podrían generar con facilidad la información antes resumida.

A	$A A A A A$	$A B A A B$	$A B B B B B$	$B B C C B$	$C C C C$	$C C C C$
	4.5		5.5		6.5	
	↑		↑		↑	
	\bar{X}_A		\bar{X}_B		\bar{X}_C	

Figura 8.13: Datos de tres muestras diferentes.

Parece evidente que los datos provienen de distribuciones con diferentes medias de población, aunque hay cierto traslape entre las muestras. Un análisis que incluya todos los datos intentaría determinar si la variabilidad entre los promedios muestrales y la variabilidad dentro de las muestras podría haber ocurrido conjuntamente *si, de hecho, las poblaciones tienen una media común*. Observe que la clave para este análisis se centra alrededor de las dos siguientes fuentes de variabilidad.

1. Variabilidad dentro de las muestras (entre observaciones en muestras distintas).
2. Variabilidad entre muestras (entre promedios muestrales).

Es evidente que si la variabilidad en 1) es considerablemente mayor que en 2), entonces habrá un traslape considerable en los datos muestrales, una señal de que los datos podrían provenir de una distribución común. En el conjunto de datos que se presenta en la

figura 8.14 se encuentra un ejemplo. Por otro lado, es muy improbable que los datos de una distribución con una media común puedan tener una variabilidad entre promedios muestrales que sea considerablemente mayor que la variabilidad dentro de las muestras.

A	B C	A C B	A C	C A B	C	A C B A	B A B A B C A C B B A B C C
						↑ ↑ ↑	
						x_A x_C x_B	

Figura 8.14: Datos que con facilidad podrían provenir de la misma población.

Las fuentes de variabilidad en 1) y 2) generan importantes cocientes de *varianzas muestrales* y los cocientes se utilizan junto con la distribución F . El procedimiento general implicado se llama **análisis de varianza**. Es interesante que en el ejemplo de la pintura aquí descrito tratamos con inferencias sobre tres medias de población pero utilizamos dos fuentes de variabilidad. No proporcionaremos detalles aquí, pero en los capítulos 13, 14 y 15 utilizaremos ampliamente el análisis de varianza en donde, por supuesto, la distribución F desempeña un papel importante.

8.8 Gráficas de cuantiles y de probabilidad

En el capítulo 1 presentamos al lector las distribuciones empíricas. El objetivo es utilizar presentaciones creativas para extraer información acerca de las propiedades de un conjunto de datos. Por ejemplo, los diagramas de tallo y hojas brindan al observador una imagen de la simetría y de otras propiedades de los datos. En este capítulo tratamos con muestras que, por supuesto, son conjuntos de datos experimentales de los que sacamos conclusiones sobre las poblaciones. A menudo, la apariencia de la muestra proporciona información sobre la distribución de la que se tomaron los datos. Por ejemplo, en el capítulo 1 ilustramos la naturaleza general de pares de muestras con gráficas de puntos que presentan una comparación relativa entre la tendencia central y la variabilidad de dos muestras.

En los capítulos siguientes con frecuencia supondremos que una distribución es normal. La información gráfica respecto a la validez de esta suposición se puede obtener a partir de presentaciones como los diagramas de tallo y hojas y los histogramas de frecuencias. Además, en esta sección presentaremos los conceptos de *gráficas de probabilidad normal* y *gráficas de cuantiles*. Estas gráficas se utilizan en estudios con diversos grados de complejidad con el principal objetivo de que las gráficas proporcionen una verificación diagnóstica sobre la suposición de que los datos provienen de una distribución normal.

Podemos caracterizar el análisis estadístico como el proceso de sacar conclusiones acerca de los sistemas en presencia de la variabilidad del sistema. Por ejemplo, el intento de un ingeniero por aprender acerca de un proceso químico a menudo es obstaculizado por la *variabilidad del proceso*. Un estudio que implica el número de artículos defectuosos en un proceso de producción con frecuencia se dificulta por la variabilidad en el método con el que se fabrican. En las secciones anteriores aprendimos acerca de las muestras y los estadísticos que expresan el centro de localización y la variabilidad en la muestra. Tales estadísticos ofrecen medidas simples, en tanto que una presentación gráfica brinda información adicional por medio de una imagen.

Un tipo de gráfica que puede ser especialmente útil para revelar la naturaleza de un conjunto de datos es la *gráfica de cuantiles*. Igual que en el caso de la gráfica de caja y extensión (véase la sección 1.6), en el que el objetivo del analista es hacer distinciones, en la gráfica de cuantiles se pueden utilizar las ideas básicas para *comparar muestras de*

datos. En los siguientes capítulos se presentarán más ejemplos del uso de las gráficas de cuantiles, en los que se analizará la inferencia estadística formal asociada con la comparación de muestras. En su momento, los estudios de caso mostrarán al lector tanto la inferencia formal como las gráficas diagnósticas para el mismo conjunto de datos.

Gráfica de cuantiles

El propósito de las gráficas de cuantiles consiste en describir, en forma de muestra, la función de distribución acumulada que se estudió en el capítulo 3.

Definición 8.6: Un **cuantil** de una muestra, $q(f)$, es un valor para el que una fracción específica f de los valores de los datos es menor que o igual a $q(f)$.

Evidentemente, un cuantil representa una estimación de una característica de una población o, más bien, la distribución teórica. La mediana de la muestra es $q(0.5)$. El percentil 75 (cuartil superior) es $q(0.75)$ y el cuartil inferior es $q(0.25)$.

Una **gráfica de cuantiles** simplemente grafica los valores de los datos en el eje vertical contra una evaluación empírica de la fracción de observaciones excedidas por los valores de los datos. Para propósitos teóricos esta fracción se calcula con

$$f_i = \frac{i - \frac{3}{8}}{n + \frac{1}{4}}$$

donde i es el orden de las observaciones cuando se ordenan de la menor a la mayor. En otras palabras, si denotamos las observaciones ordenadas como

$$y_{(1)} \leq y_{(2)} \leq y_{(3)} \leq \dots \leq y_{(n-1)} \leq y_{(n)},$$

entonces la gráfica de cuantiles describe una gráfica de $y_{(i)}$ contra f_i . En la figura 8.15 se presenta la gráfica de cuantiles para las asas de las latas de pintura analizadas con anterioridad.

A diferencia de la gráfica de caja y extensión, la gráfica de cuantiles realmente muestra todas las observaciones. Todos los cuantiles, incluidos la mediana y los cuantiles superior e inferior, se pueden aproximar de forma visual. Por ejemplo, observamos fácilmente una mediana de 35 y un cuartil superior de alrededor de 36. Las agrupaciones relativamente grandes en torno a valores específicos se indican por pendientes cercanas a cero; mientras que los datos escasos en ciertas áreas producen pendientes más abruptas. La figura 8.15 describe la dispersión de datos de los valores 28 a 30, pero una densidad relativamente alta de 36 a 38. En los capítulos 9 y 10 proseguimos con las gráficas de cuantiles mediante la ilustración de formas útiles en que es posible comparar distintas muestras.

Debería ser muy evidente para el lector que detectar si un conjunto de datos proviene o no de una distribución normal puede ser una herramienta importante para el analista de datos. Como antes indicamos en esta sección, a menudo suponemos que la totalidad o subconjuntos de las observaciones en un conjunto de datos son realizaciones de variables aleatorias normales independientes idénticamente distribuidas. Una vez más, la gráfica de diagnóstico a menudo se agrega a (con fines de presentación) una *prueba de bondad del ajuste* formal de los datos. Las pruebas de bondad del ajuste se estudiarán en el capítulo 10. Los lectores de un artículo o informe científico suelen considerar la información de diagnóstico mucho más clara, menos árida y quizá menos aburrida que un análisis formal.

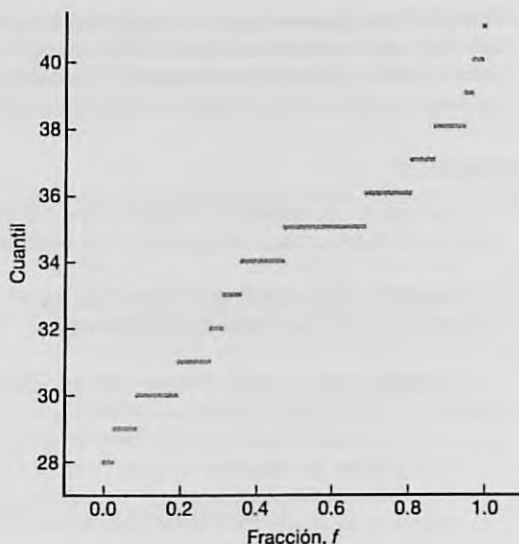


Figura 8.15: Gráfica de cuantiles para los datos de la pintura.

En los capítulos siguientes (del 9 al 13) nos enfocaremos nuevamente en los métodos de detección de desviaciones de la normalidad como un agregado de la inferencia estadística formal. Las gráficas de cuantiles son útiles para detectar los tipos de distribución. En la elaboración de modelos y en el diseño de experimentos también hay situaciones en que se utilizan las gráficas para detectar **términos o efectos del modelo** que están activos. En otras situaciones se utilizan para determinar si las suposiciones subyacentes que el científico o el ingeniero hicieron en la construcción del modelo son o no razonables. En los capítulos 11, 12 y 13 se incluyen muchos ejemplos con ilustraciones. La siguiente subsección brinda un análisis y un ejemplo de una gráfica de diagnóstico denominada *gráfica de cuantiles-cuantiles normales*.

Gráfica de cuantiles-cuantiles normales

La gráfica de cuantiles-cuantiles normales aprovecha lo que se conoce sobre los cuantiles de la distribución normal. La metodología incluye una gráfica de los cuantiles empíricos recién analizados, contra el cuantil correspondiente de la distribución normal. Ahora, la expresión para un cuantil de una variable aleatoria $N(\mu, \sigma)$ es muy complicada. Sin embargo, una buena aproximación es dada por

$$q_{\mu, \sigma}(f) = \mu + \sigma\{4.91[f^{0.14} - (1-f)^{0.14}]\}.$$

La expresión entre las llaves (el múltiplo de σ) es la aproximación para el cuantil correspondiente para la variable aleatoria $N(0, 1)$, es decir,

$$q_{0,1}(f) = 4.91[f^{0.14} - (1-f)^{0.14}].$$

Definición 8.7: La **gráfica de cuantiles-cuantiles normales** es una gráfica de $y_{(i)}$ (observaciones ordenadas) contra $q_{0,1}(f_i)$, donde $f_i = \frac{i - \frac{1}{2}}{n + \frac{1}{2}}$.

Una relación cercana a una línea recta sugiere que los datos provienen de una distribución normal. La intersección en el eje vertical es una estimación de la media de la población μ y la pendiente es una estimación de la desviación estándar σ . La figura 8.16 presenta una gráfica de cuantiles-cuantiles normales para los datos de las latas de pintura.

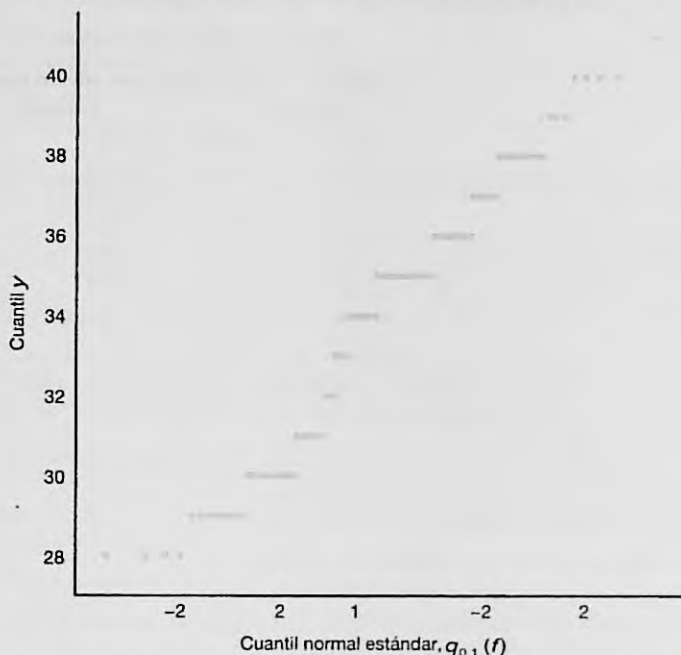


Figura 8.16: Gráfica de cuantiles-cuantiles normales para los datos de la pintura.

Graficación de la probabilidad normal

Observe cómo la desviación de la normalidad se vuelve evidente gracias a la apariencia de la gráfica. La asimetría que exhiben los datos produce cambios en la pendiente.

Las ideas para graficar la probabilidad se manifiestan en versiones diferentes de la gráfica de cuantiles-cuantiles normales que se presentó aquí. Por ejemplo, se ha puesto mucha atención a la llamada **gráfica de probabilidad normal**, en la que f se grafica contra los valores de los datos ordenados en un papel especial y la escala utilizada da como resultado una línea recta. Además, una gráfica alternativa utiliza los valores esperados de las observaciones clasificadas para la distribución normal y dibuja las observaciones clasificadas contra su valor esperado, bajo el supuesto de datos de $N(\mu, \sigma)$. Una vez más, la línea recta es el criterio gráfico que se emplea. Continuamos sugiriendo que basarse en los métodos analíticos gráficos que se describen en esta sección ayudará a comprender los métodos formales que permiten distinguir muestras diferentes de datos.

Ejemplo 8.12: Considere los datos del ejercicio 10.41 en la página 358 del capítulo 10. En el estudio “Retención de nutrientes y respuesta de comunidades de macroinvertebrados ante la presión de aguas residuales en un ecosistema fluvial”, que se llevó a cabo en el departamento de zoología del Virginia Polytechnic Institute y la universidad estatal, se recabaron datos sobre mediciones de densidad (número de organismos por metro cuadrado) en dos diferentes estaciones colectoras. En el capítulo 10 se dan detalles con respecto a los métodos analíticos de comparación de muestras para determinar si ambas provienen de la misma distribución $N(\mu, \sigma)$. Los datos se presentan en la tabla 8.1.

Tabla 8.1: Datos para el ejemplo 8.12

Número de organismos por metro cuadrado			
Estación 1		Estación 2	
5,030	4,980	2,800	2,810
13,700	11,910	4,670	1,330
10,730	8,130	6,890	3,320
11,400	26,850	7,720	1,230
860	17,660	7,030	2,130
2,200	22,800	7,330	2,190
4,250	1,130		
15,040	1,690		

Dibuje una gráfica de cuantiles-cuantiles normales y saque conclusiones con respecto a si es razonable o no suponer que las dos muestras provienen de la misma distribución $n(x; \mu, \sigma)$.

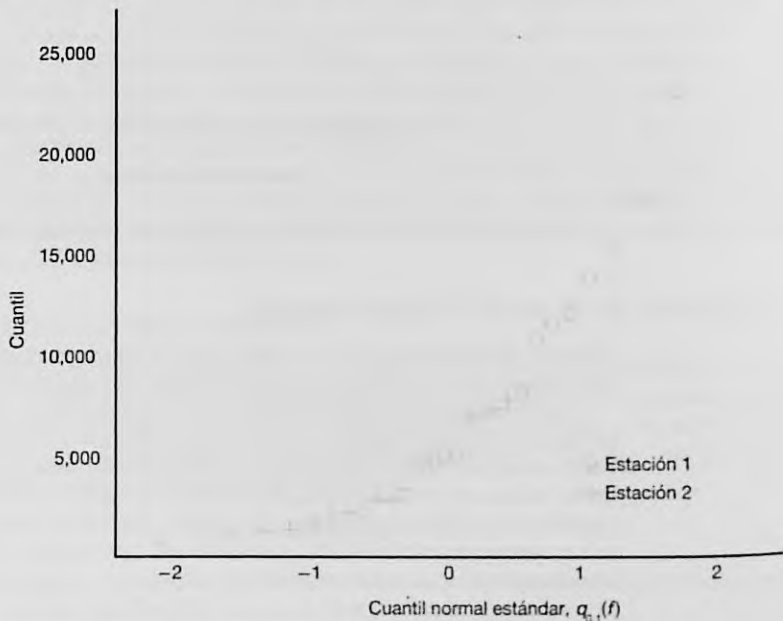


Figura 8.17: Gráfica de cuantiles-cuantiles normales para los datos de densidad del ejemplo 8.12.

Solución: La figura 8.17 muestra la gráfica de cuantiles-cuantiles normales para las mediciones de densidad. La gráfica se aleja mucho de una sola línea recta. De hecho, los datos de la estación 1 reflejan pocos valores en la cola inferior de la distribución y varios en la cola superior. El "agrupamiento" de observaciones hace que parezca improbable que las dos muestras provengan de una distribución común $N(\mu, \sigma)$. ■

Aunque hemos concentrado nuestra explicación y ejemplo en las gráficas de probabilidad para distribuciones normales, podemos enfocarnos en cualquier distribución. Tan sólo necesitaríamos calcular cantidades de forma analítica para la distribución teórica en cuestión.

Ejercicios

8.37 Para una distribución chi cuadrada calcule

- $\chi_{0.025}^2$ cuando $\nu = 15$;
- $\chi_{0.01}^2$ cuando $\nu = 7$;
- $\chi_{0.05}^2$ cuando $\nu = 24$.

8.38 Para una distribución chi cuadrada, calcule

- $\chi_{0.005}^2$ cuando $\nu = 5$;
- $\chi_{0.05}^2$ cuando $\nu = 19$;
- $\chi_{0.01}^2$ cuando $\nu = 12$.

8.39 Para una distribución chi cuadrada calcule χ_{α}^2 , tal que

- $P(X^2 > \chi_{\alpha}^2) = 0.99$ cuando $\nu = 4$;
- $P(X^2 > \chi_{\alpha}^2) = 0.025$ cuando $\nu = 19$;
- $P(37.652 < X^2 < \chi_{\alpha}^2) = 0.045$ cuando $\nu = 25$.

8.40 Para una distribución chi cuadrada calcule χ_{α}^2 , tal que

- $P(X^2 > \chi_{\alpha}^2) = 0.01$ cuando $\nu = 21$;
- $P(X^2 < \chi_{\alpha}^2) = 0.95$ cuando $\nu = 6$;
- $P(\chi_{\alpha}^2 < X^2 < 23.209) = 0.015$ cuando $\nu = 10$.

8.41 Suponga que las varianzas muestrales son mediciones continuas. Calcule la probabilidad de que una muestra aleatoria de 25 observaciones, de una población normal con varianza $\sigma^2 = 6$, tenga una varianza muestral S^2

- mayor que 9.1;
- entre 3.462 y 10.745.

8.42 Las calificaciones de un examen de colocación que se aplicó a estudiantes de primer año de una universidad durante los últimos cinco años tienen una distribución aproximadamente normal con una media $\mu = 74$ y una varianza $\sigma^2 = 8$. ¿Seguiría considerando que $\sigma^2 = 8$ es un valor válido de la varianza si una muestra aleatoria de 20 estudiantes, a los que se les aplica el

examen de colocación este año, obtienen un valor de $s^2 = 20$?

8.43 Demuestre que la varianza de S^2 para muestras aleatorias de tamaño n de una población normal disminuye a medida que aumenta n . [Sugerencia: primero calcule la varianza de $(n-1)S^2/\sigma^2$].

8.44 a) Calcule $t_{0.025}$ cuando $\nu = 14$.

b) Calcule $-t_{0.10}$ cuando $\nu = 10$.

c) Calcule $t_{0.995}$ cuando $\nu = 7$.

8.45 a) Calcule $P(T < 2.365)$ cuando $\nu = 7$.

b) Calcule $P(T > 1.318)$ cuando $\nu = 24$.

c) Calcule $P(-1.356 < T < 2.179)$ cuando $\nu = 12$.

d) Calcule $P(T > -2.567)$ cuando $\nu = 17$.

8.46 a) Calcule $P(-t_{0.005} < T < t_{0.01})$ para $\nu = 20$.

b) Calcule $P(T > -t_{0.025})$.

8.47 Dada una muestra aleatoria de tamaño 24 de una distribución normal, calcule k tal que

a) $P(-2.069 < T < k) = 0.965$;

b) $P(k < T < 2.807) = 0.095$;

c) $P(-k < T < k) = 0.90$.

8.48 Una empresa que fabrica juguetes electrónicos afirma que las baterías que utiliza en sus productos duran un promedio de 30 horas. Para mantener este promedio se prueban 16 baterías cada mes. Si el valor t calculado cae entre $-t_{0.025}$ y $t_{0.025}$, la empresa queda satisfecha con su afirmación. ¿Qué conclusiones debería sacar la empresa a partir de una muestra que tiene una media de $\bar{x} = 27.5$ horas y una desviación estándar de $s = 5$ horas? Suponga que la distribución de las duraciones de las baterías es aproximadamente normal.

8.49 Una población normal con varianza desconocida tiene una media de 20. ¿Es posible obtener una muestra aleatoria de tamaño 9 de esta población con una media de 24 y una desviación estándar de 4.1? Si no fuera posible, ¿a qué conclusión llegaría?

8.50 Un fabricante de cierta marca de barras de cereal con bajo contenido de grasa afirma que el contenido promedio de grasa saturada en éstas es de 0.5 gramos. En una muestra aleatoria de 8 barras de cereal de esta marca se encontró que su contenido de grasa saturada era de 0.6, 0.7, 0.7, 0.3, 0.4, 0.5, 0.4 y 0.2. ¿Estaría de acuerdo con tal afirmación? Suponga una distribución normal.

8.51 Para una distribución F calcule:

- $f_{0.05}$ con $v_1 = 7$ y $v_2 = 15$;
- $f_{0.05}$ con $v_1 = 15$ y $v_2 = 7$;
- $f_{0.01}$ con $v_1 = 24$ y $v_2 = 19$;
- $f_{0.95}$ con $v_1 = 19$ y $v_2 = 24$;
- $f_{0.99}$ con $v_1 = 28$ y $v_2 = 12$.

8.52 Se aplican pruebas a 10 cables conductores soldados a un dispositivo semiconductor con el fin de determinar su resistencia a la tracción. Las pruebas demostraron que para romper la unión se requieren las libras de fuerza que se listan a continuación:

19.8	12.7	13.2	16.9	10.6
18.8	11.1	14.3	17.0	12.5

Otro conjunto de 8 cables conductores que forman un dispositivo se encapsuló y se probó para determinar si el encapsulado aumentaba la resistencia a la tracción. Las pruebas dieron los siguientes resultados:

24.9	22.8	23.6	22.1	20.4	21.6	21.8	22.5
------	------	------	------	------	------	------	------

Comente acerca de la evidencia disponible respecto a la igualdad de las dos varianzas de población.

8.53 Considere las siguientes mediciones de la capa-

cidad de producción de calor del carbón producido por dos minas (en millones de calorías por tonelada):

Mina 1:	8260	8130	8350	8070	8340	
Mina 2:	7950	7890	7900	8140	7920	7840

¿Se puede concluir que las dos varianzas de población son iguales?

8.54 Dibuje una gráfica de cuantiles con los siguientes datos, que representan la vida, en horas, de cincuenta lámparas incandescentes esmeriladas de 40 watts y 110 voltios, tomados de pruebas de vida forzadas:

919	1196	785	1126	936	918
1156	920	948	1067	1092	1162
1170	929	950	905	972	1035
1045	855	1195	1195	1340	1122
938	970	1237	956	1102	1157
978	832	1009	1157	1151	1009
765	958	902	1022	1333	811
1217	1085	896	958	1311	1037
702	923				

8.55 Dibuje una gráfica de cuantiles-cuantiles normales con los siguientes datos, que representan los diámetros de 36 cabezas de remache en 1/100 de una pulgada:

6.72	6.77	6.82	6.70	6.78	6.70	6.62
6.75	6.66	6.66	6.64	6.76	6.73	6.80
6.72	6.76	6.76	6.68	6.66	6.62	6.72
6.76	6.70	6.78	6.76	6.67	6.70	6.72
6.74	6.81	6.79	6.78	6.66	6.76	6.76
6.72						

Ejercicios de repaso

8.56 Considere los datos que se presentan en el ejercicio 1.20 de la página 31. Dibuje una gráfica de caja y extensión, y comente acerca de la naturaleza de la muestra. Calcule la media muestral y la desviación estándar de la muestra.

8.57 Si X_1, X_2, \dots, X_n son variables aleatorias independientes que tienen distribuciones exponenciales idénticas con parámetro θ , demuestre que la función de densidad de la variable aleatoria $Y = X_1 + X_2 + \dots + X_n$ es la de una distribución gamma con parámetros $\alpha = n$ y $\beta = \theta$.

8.58 Al probar el monóxido de carbono que contiene cierta marca de cigarrillos, los datos que se obtuvieron, en miligramos por cigarrillo, se codificaron restando 12 a cada observación. Utilice los resultados del ejercicio 8.14 de la página 231 para calcular la desviación estándar del contenido de monóxido de carbono de una muestra aleatoria de 15 cigarrillos de esta marca, si las mediciones codificadas son 3.8, -0.9, 5.4, 4.5, 5.2, 5.6, -0.1, -0.3, -1.7, 5.7, 3.3, 4.4, -0.5 y 1.9.

8.59 Si S_1^2 y S_2^2 representan las varianzas de muestras aleatorias independientes de tamaños $n_1 = 8$ y $n_2 = 12$, tomadas de poblaciones normales con varianzas iguales, calcule $P(S_1^2 / S_2^2 < 4.89)$.

8.60 Una muestra aleatoria de 5 presidentes de bancos indicó sueldos anuales de \$395,000, \$521,000, \$483,000, \$479,000 y \$510,000. Calcule la varianza de este conjunto.

8.61 Si el número de huracanes que azotan cierta área del este de Estados Unidos cada año es una variable aleatoria que tiene una distribución de Poisson con $\mu = 6$, calcule la probabilidad de que esta área sea azotada por

- exactamente 15 huracanes en 2 años;
- a lo sumo 9 huracanes en 2 años.

8.62 Una empresa de taxis prueba una muestra aleatoria de 10 neumáticos radiales con bandas tensoras de acero de cierta marca y registra los siguientes desgastes de la banda: 48,000, 53,000, 45,000, 61,000, 59,000, 56,000, 63,000, 49,000, 53,000 y 54,000 kilómetros.

Utilice los resultados del ejercicio 8.14 de la página 231 para calcular la desviación estándar de este conjunto de datos dividiendo primero cada observación entre 1000 y después restando 55 al resultado.

8.63 Considere los datos del ejercicio 1.19 de la página 31. Dibuje una gráfica de caja y extensión. Comente y calcule la media muestral y la desviación estándar muestral.

8.64 Si S_1^2 y S_2^2 representan las varianzas de muestras aleatorias independientes de tamaños $n_1 = 25$ y $n_2 = 31$, tomadas de poblaciones normales con varianzas $\sigma_1^2 = 10$ y $\sigma_2^2 = 15$, respectivamente, calcule

$$P(S_1^2/S_2^2 > 1.26).$$

8.65 Considere el ejemplo 1.5 de la página 25. Comente acerca de cualquier valor extremo.

8.66 Considere el ejercicio de repaso 8.56. Comente acerca de cualquier valor extremo en los datos.

8.67 La resistencia a la rotura X de cierto remache que se utiliza en el motor de una máquina tiene una media de 5000 psi y una desviación estándar de 400 psi. Se toma una muestra aleatoria de 36 remaches. Considere la distribución de \bar{X} , la media muestral de la resistencia a la rotura.

a) ¿Cuál es la probabilidad de que la media de la muestra caiga entre 4800 psi y 5200 psi?

b) ¿Qué muestra n sería necesaria para tener

$$P(4900 < \bar{X} < 5100) = 0.99?$$

8.68 Considere la situación del ejercicio de repaso 8.62. Si la población de la cual se tomó la muestra tiene una media poblacional $\mu = 53,000$ kilómetros, ¿esta información de la muestra parece apoyar esa afirmación? En su respuesta calcule

$$t = \frac{\bar{x} - 53,000}{s/\sqrt{10}}$$

y determine, consultando la tabla A.4 (con 9 g.l.), si el valor t calculado es razonable o si parece ser un suceso raro.

8.69 Se consideran dos propulsores de combustible sólido distintos, el tipo A y el tipo B , para una actividad del programa espacial. Las velocidades de combustión en el propulsor son fundamentales. Se toman muestras aleatorias de 20 especímenes de los dos propulsores con medias muestrales de 20.5 cm/s para el propulsor A y de 24.50 cm/s para el propulsor B . Por lo general se supone que la variabilidad en la velocidad de combustión es casi igual para los dos propulsores y que es determinada por una desviación estándar de población de 5 cm/s. Suponga que la velocidad de combustión

para cada propulsor es aproximadamente normal, por lo cual se debería utilizar el teorema del límite central. Nada se sabe acerca de las medias poblacionales de las dos velocidades de combustión y se espera que este experimento revele algo sobre ellas.

a) Si, de hecho, $\mu_A = \mu_B$, ¿cuál será $P(\bar{X}_B - \bar{X}_A \geq 4.0)$?

b) Utilice lo que respondió en el inciso a) para dar luz sobre la validez de la proposición $\mu_A = \mu_B$.

8.70 La concentración de un ingrediente activo en el producto de una reacción química es fuertemente influido por el catalizador que se usa en la reacción. Se considera que cuando se utiliza el catalizador A la concentración media de la población excede el 65%. Se sabe que la desviación estándar es $\sigma = 5\%$. Una muestra de productos tomada de 30 experimentos independientes proporciona la concentración promedio de $\bar{x}_A = 64.5\%$.

a) ¿Esta información muestral, con una concentración promedio de $\bar{x}_A = 64.5\%$, ofrece información inquietante de que quizá μ_A no sea el 65% sino menos que ese porcentaje? Respalde su respuesta con una aseveración de probabilidad.

b) Suponga que se realiza un experimento similar utilizando otro catalizador, el B . Se supone que la desviación estándar σ sigue siendo 5% y \bar{x}_B resulta ser 70%. Comente si la información muestral del catalizador B sugiere con certeza que μ_B es en realidad mayor que μ_A . Respalde su respuesta calculando

$$P(\bar{X}_B - \bar{X}_A \geq 5.5 \mid \mu_B = \mu_A).$$

c) En el caso de que $\mu_A = \mu_B = 65\%$, determine la distribución aproximada de las siguientes cantidades (con la media y la varianza de cada una). Utilice el teorema del límite central.

- i) \bar{X}_B ;
- ii) $\bar{X}_A - \bar{X}_B$;
- iii) $\frac{\bar{X}_A - \bar{X}_B}{\sigma\sqrt{2/30}}$.

8.71 Con la información del ejercicio de repaso 8.70 calcule (suponiendo $\mu_B = 65\%$) $P(\bar{X}_B \geq 70)$.

8.72 Dada una variable aleatoria normal X con media 20 y varianza 9, y una muestra aleatoria de tamaño n tomada de la distribución, ¿qué tamaño de la muestra n se necesita para que

$$P(19.9 \leq \bar{X} \leq 20.1) = 0.95?$$

8.73 En el capítulo 9 se estudiará con detenimiento el concepto de estimación de parámetros. Suponga que X es una variable aleatoria con media μ y varianza $\sigma^2 = 1.0$. Además, suponga que se toma una muestra aleato-

ria de tamaño n y que \bar{x} se utiliza como un *estimado* de μ . Cuando se toman los datos y se mide la media de la muestra, deseamos que ésta esté dentro de 0.05 unidades de la media real con una probabilidad de 0.99. Es decir, aquí queremos que haya muchas posibilidades de que la \bar{x} calculada de la muestra esté “muy cerca de” la media de población (¡dondequiera que ésta se encuentre!), de manera que deseamos

$$P(|\bar{X} - \mu| > 0.05) = 0.99.$$

¿Qué tamaño de muestra se requiere?

8.74 Suponga que se utiliza una máquina para llenar envases de cartón con un líquido. La especificación que es estrictamente indispensable para el llenado de la máquina es 9 ± 1.5 onzas. El proveedor considera que cualquier envase de cartón que no cumpla con tales límites de peso en el llenado está defectuoso. Se espera que al menos 99% de los envases de cartón cumplan con la especificación. En el caso de que $\mu = 9$ y $\sigma = 1$, ¿qué proporción de envases de cartón del proceso están defectuosos? Si se hacen cambios para reducir la variabilidad, ¿cuánto se tiene que reducir σ para que haya 0.99 de probabilidades de cumplir con la especificación? Suponga una distribución normal para el peso.

8.75 Considere la situación del ejercicio de repaso 8.74. Suponga que se hace un gran esfuerzo para “estrechar” la variabilidad del sistema. Después de eso se toma una muestra aleatoria de tamaño 40 de la nueva

línea de ensamble y se obtiene que la varianza de la muestra es $s^2 = 0.188$ onzas². ¿Tenemos evidencia numérica sólida de que σ^2 se redujo a menos de 1.0? Considere la probabilidad

$$P(S^2 \leq 0.188 \mid \sigma^2 = 1.0),$$

y dé una conclusión.

8.76 Proyecto de grupo: Divida al grupo en equipos de cuatro estudiantes. Cada equipo deberá ir al gimnasio de la universidad o a un gimnasio local y preguntar a cada persona que cruce el umbral cuánto mide en pulgadas. Después, cada equipo dividirá los datos de las estaturas por género y trabajará en conjunto para realizar las actividades que se indican a continuación.

- Dibujen una gráfica de cuantiles-cuantiles normal con los datos. Si usan la gráfica como base, ¿les parecería que los datos tienen una distribución normal?
- Utilicen la varianza muestral como un estimado de la varianza real para cada género. Supongan que la estatura media de la población de los hombres es realmente tres pulgadas más grande que la de las mujeres. ¿Cuál es la probabilidad de que la estatura promedio de los hombres sea 4 pulgadas más grande que la de las mujeres en su muestra?
- ¿Qué factores podrían provocar que estos resultados sean engañosos?

8.9 Posibles riesgos y errores conceptuales. Relación con el material de otros capítulos

El teorema del límite central es una de las más poderosas herramientas de la estadística, y aunque este capítulo es relativamente breve, contiene gran cantidad de información fundamental acerca de las herramientas que se utilizarán en el resto del libro.

El concepto de distribución muestral es una de las ideas fundamentales más importantes de la estadística y, en este momento de su entrenamiento, el estudiante debería entenderlo con claridad antes de continuar con los siguientes capítulos, en los cuales se continuarán utilizando ampliamente las distribuciones muestrales. Suponga que se quiere utilizar el estadístico \bar{X} para hacer inferencias acerca de la media de la población μ , lo cual se hace utilizando el valor observado \bar{x} de una sola muestra de tamaño n . Luego, cualquier inferencia deberá hacerse tomando en cuenta no sólo el valor único, sino también la estructura teórica o la **distribución de todos los valores \bar{x} que se podrían observar a partir de las muestras de tamaño n** . Como resultado de lo anterior surge el concepto de *distribución muestral*, que es la base del teorema del límite central. Las distribuciones t , χ^2 y F también se utilizan en el contexto de las distribuciones muestrales. Por ejemplo, la distribución t , que se ilustra en la figura 8.8, representa la estructura que ocurre si se forman todos los valores de $\frac{\bar{x} - \mu}{s/\sqrt{n}}$, donde \bar{x} y s se toman de las

muestras de tamaño n de una distribución $n(x; \mu, \sigma)$. Se pueden hacer comentarios similares en relación con χ^2 y F , y el lector no debería olvidar que la información muestral que conforma los estadísticos para todas estas distribuciones es la normal. Por lo tanto, se podría afirmar que **donde haya una t , F o χ^2 la fuente era una muestra de una distribución normal.**

Podría parecer que las tres distribuciones antes descritas se presentaron de una forma bastante aislada, sin indicar a qué se refieren. Sin embargo, aparecerán en la resolución de problemas prácticos a lo largo del texto.

Ahora bien, hay tres cuestiones que se deben tener presentes para evitar que haya confusión respecto a estas distribuciones muestrales fundamentales:

- i) No se puede usar el teorema del límite central a menos que se conozca σ . Para usar el teorema del límite central cuando no se conoce σ se debe reemplazar con s , la desviación estándar de la muestra.
- ii) El estadístico T no es un resultado del teorema del límite central y x_1, x_2, \dots, x_n deben provenir de una distribución $n(x; \mu, \sigma)$ para que $\frac{\bar{x} - \mu}{s/\sqrt{n}}$ sea una distribución t ; por supuesto, s es tan sólo una estimación de σ .
- iii) Aunque el concepto de **grados de libertad** es nuevo en este punto, debería ser muy intuitivo, ya que es razonable que la naturaleza de la distribución de S y también t deban depender de la cantidad de información en la muestra x_1, x_2, \dots, x_n .

Capítulo 9

Problemas de estimación de una y dos muestras

9.1 Introducción

En los capítulos anteriores destacamos las propiedades del muestreo de la media y de la varianza muestrales. También destacamos las representaciones de datos en varias formas. El propósito de estas presentaciones es establecer las bases que permitan a los estadísticos sacar conclusiones acerca de los parámetros de poblaciones tomadas de datos experimentales. Por ejemplo, el teorema del límite central brinda información sobre la distribución de la media muestral \bar{X} . La distribución incluye la media de la población μ . Por consiguiente, cualesquiera conclusiones respecto a μ , extraídas de un promedio muestral observado, deben depender de lo que se sabe acerca de su distribución muestral. Se podría decir algo similar en lo que se refiere a S^2 y σ^2 . Como es evidente, es muy probable que cualquier conclusión que saquemos acerca de la varianza de una distribución normal implique la distribución muestral de S^2 .

En este capítulo comenzaremos por presentar de manera formal el propósito de la inferencia estadística. Continuaremos con el análisis del problema de la **estimación de los parámetros de la población**. Restringiremos nuestros desarrollos formales de los procedimientos de estimación específicos a problemas que impliquen una y dos muestras.

9.2 Inferencia estadística

En el capítulo 1 presentamos la filosofía general de la inferencia estadística formal. La **inferencia estadística** consta de los métodos mediante los cuales se hacen inferencias o generalizaciones acerca de una población. La tendencia actual es distinguir entre el **método clásico** de estimación de un parámetro de la población, donde las inferencias se basan estrictamente en información obtenida de una muestra aleatoria seleccionada de la población, y el **método bayesiano**, el cual utiliza el conocimiento subjetivo que ya se posee sobre la distribución de probabilidad de los parámetros desconocidos junto con la información que proporcionan los datos de la muestra. En la mayor parte de este capítulo utilizaremos los métodos clásicos para estimar los parámetros de la población desconocidos, como la media, la proporción y la varianza, mediante el cálculo de estadísticos de muestras aleatorias y la aplicación de la teoría de las distribuciones muestrales, gran

parte de lo cual se estudió en el capítulo 8. La estimación bayesiana se analizará en el capítulo 18.

La inferencia estadística se puede dividir en dos áreas principales: **estimación y pruebas de hipótesis**. Trataremos estas dos áreas por separado: en este capítulo veremos la teoría y las aplicaciones de la estimación, y en el capítulo 10 revisaremos la prueba de hipótesis. Para distinguir claramente un área de la otra, considere los siguientes ejemplos. Un candidato a un cargo público podría estar interesado en estimar la verdadera proporción de votantes que lo favorecerán mediante la obtención de las opiniones de una muestra aleatoria de 100 de ellos. La parte de votantes en la muestra que favorecerán al candidato se podría utilizar como un estimado de la verdadera proporción en la población de votantes. El conocimiento de la distribución muestral de una proporción nos permite establecer el grado de exactitud de tal estimado. Este problema cae en el área de la estimación.

Considere ahora el caso de alguien a quien le interesa averiguar si la marca A de cera para piso es más resistente al desgaste que la marca B . Se podría plantear la hipótesis de que la marca A es mejor que la marca B y, después de la prueba adecuada, aceptar o rechazar dicha hipótesis. En este ejemplo no intentamos estimar un parámetro, sino llegar a una decisión correcta acerca de una hipótesis planteada previamente. Una vez más, dependemos de la teoría del muestreo y de utilizar datos que nos proporcionen alguna medida del grado de exactitud de nuestra decisión.

9.3 Métodos de estimación clásicos

La **estimación puntual** de algún parámetro de la población θ es un solo valor $\hat{\theta}$ de un estadístico $\hat{\Theta}$. Por ejemplo, el valor \bar{x} del estadístico \bar{X} , que se calcula a partir de una muestra de tamaño n , es una estimación puntual del parámetro de la población μ . De manera similar, $\hat{p} = x/n$ es una estimación puntual de la verdadera proporción p para un experimento binomial.

No se espera que un estimador logre estimar el parámetro de la población sin error. No se espera que \bar{X} estime μ con exactitud, lo que en realidad se espera es que no esté muy alejada. Para una muestra específica, la manera en que se podría obtener un estimado más cercano de μ es utilizando la mediana de la muestra \bar{X} como estimador. Considere, por ejemplo, una muestra que consta de los valores 2, 5 y 11 de una población cuya media es 4, la cual, supuestamente, se desconoce. Podríamos estimar μ para que sea $\bar{x} = 6$ usando la media muestral como nuestro estimado, o bien, $\bar{x} = 5$ utilizando la mediana muestral. En este caso el estimador \bar{X} produce una estimación más cercana al parámetro verdadero que la que produce el estimador X^- . Por otro lado, si nuestra muestra aleatoria contiene los valores 2, 6 y 7, entonces $\bar{x} = 5$ y $\bar{x} = 6$, de manera que el mejor estimador es \bar{X} . Cuando no conocemos el valor real de μ , tenemos que comenzar por decidir qué estimador utilizaremos, si X^- o \bar{X} .

Estimador insesgado

¿Cuáles son las propiedades que una “buena” función de decisión debería tener para poder influir en nuestra elección de un estimador en vez de otro? Sea $\hat{\Theta}$ un estimador cuyo valor $\hat{\theta}$ es una estimación puntual de algún parámetro de la población desconocido θ . Sin duda deseáramos que la distribución muestral de $\hat{\Theta}$ tuviera una media igual al parámetro estimado. Al estimador que tuviera esta propiedad se le llamaría **estimador insesgado**.

Definición 9.1: Se dice que un estadístico $\hat{\Theta}$ es un **estimador insesgado** del parámetro θ si

$$\mu_{\hat{\Theta}} = E(\hat{\Theta}) = \theta.$$

Ejemplo 9.1: Demuestre que S^2 es un estimador insesgado del parámetro σ^2 .

Solución: En la sección 8.5, en la página 244, demostramos que

$$\sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2.$$

Entonces,

$$\begin{aligned} E(S^2) &= E\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n E(X_i - \mu)^2 - nE(\bar{X} - \mu)^2 \right] = \frac{1}{n-1} \left(\sum_{i=1}^n \sigma_{X_i}^2 - n\sigma_{\bar{X}}^2 \right). \end{aligned}$$

Sin embargo,

$$\sigma_{X_i}^2 = \sigma^2, \text{ para } i = 1, 2, \dots, n, \text{ y } \sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}.$$

Por lo tanto,

$$E(S^2) = \frac{1}{n-1} \left(n\sigma^2 - n\frac{\sigma^2}{n} \right) = \sigma^2. \quad \blacksquare$$

Aunque S^2 es un estimador insesgado de σ^2 , S , por otro lado, suele ser un estimador sesgado de σ , un sesgo que en el caso de muestras grandes se vuelve insignificante. Este ejemplo ilustra **por qué dividimos entre $n-1$** en vez de entre n cuando estimamos la varianza.

Varianza de un estimador puntual

Si $\hat{\Theta}_1$ y $\hat{\Theta}_2$ son dos estimadores insesgados del mismo parámetro de la población θ , deseamos elegir el estimador cuya distribución muestral tenga la menor varianza. Por lo tanto, si $\sigma_{\hat{\Theta}_1}^2 < \sigma_{\hat{\Theta}_2}^2$, decimos que $\hat{\Theta}_1$ es un **estimador más eficaz** de θ que $\hat{\Theta}_2$.

Definición 9.2: Si consideramos todos los posibles estimadores insesgados de algún parámetro θ , al que tiene la menor varianza lo llamamos **estimador más eficaz** de θ .

En la figura 9.1 se ilustran las distribuciones muestrales de tres estimadores diferentes $\hat{\Theta}_1$, $\hat{\Theta}_2$ y $\hat{\Theta}_3$, todos para θ . Es evidente que sólo $\hat{\Theta}_1$ y $\hat{\Theta}_2$ no son sesgados, ya que sus distribuciones están centradas en θ . El estimador $\hat{\Theta}_1$ tiene una varianza menor que $\hat{\Theta}_2$, por lo tanto, es más eficaz. En consecuencia, el estimador de θ que elegiríamos, entre los tres que estamos considerando, sería $\hat{\Theta}_1$.

Para poblaciones normales se puede demostrar que tanto \bar{X} como \bar{x} son estimadores insesgados de la media de la población μ , pero la varianza de \bar{X} es más pequeña que la varianza de \bar{x} . Por consiguiente, los estimados \bar{x} y \bar{X} serán, en promedio, iguales a

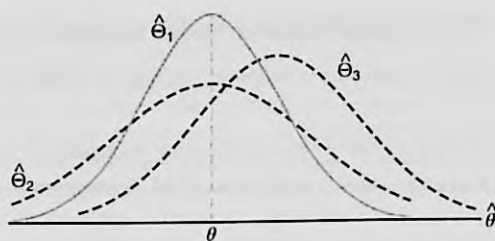


Figura 9.1: Distribuciones muestrales de diferentes estimadores de θ .

la media de la población μ , aunque podría ser que \bar{x} esté más cerca de μ para una muestra dada y, por lo tanto, que \bar{X} sea más eficaz que \bar{X} .

Estimación por intervalo

Podría ser que ni el estimador insesgado más eficaz estime con exactitud el parámetro de la población. Es cierto que la exactitud de la estimación aumenta cuando las muestras son grandes; pero incluso así no tenemos razones para esperar que una **estimación puntual** de una muestra dada sea exactamente igual al parámetro de la población que se supone debe estimar. Hay muchas situaciones en que es preferible determinar un intervalo dentro del cual esperaríamos encontrar el valor del parámetro. Tal intervalo se conoce como **estimación por intervalo**.

Una estimación por intervalo de un parámetro de la población θ es un intervalo de la forma $\hat{\theta}_L < \theta < \hat{\theta}_U$, donde $\hat{\theta}_L$ y $\hat{\theta}_U$ dependen del valor del estadístico $\hat{\Theta}$ para una muestra específica, y también de la distribución de muestreo de $\hat{\Theta}$. Por ejemplo, una muestra aleatoria de calificaciones verbales de la prueba SAT para estudiantes universitarios de primer año produciría un intervalo de 530 a 550, dentro del cual esperamos encontrar el promedio verdadero de todas las calificaciones verbales de la prueba SAT para ese grupo. Los valores de los puntos extremos, 530 y 550, dependerán de la media muestral calculada \bar{x} y de la distribución de muestreo de \bar{X} . A medida que aumenta el tamaño de la muestra, sabemos que $\sigma_{\bar{X}}^2 = \sigma^2/n$ disminuye y, en consecuencia, cabe la posibilidad de que nuestra estimación se acerque más al parámetro μ , lo cual daría como resultado un intervalo más corto. De esta manera, el intervalo de la estimación indica, por su longitud, la precisión de la estimación puntual. Un ingeniero obtendrá información acerca de la proporción de la población de artículos defectuosos tomando una muestra y calculando la *proporción muestral defectuosa*, sin embargo, una estimación por intervalo podría ser más informativa.

Interpretación de las estimaciones por intervalo

Como muestras distintas suelen producir valores diferentes de $\hat{\Theta}$ y, por lo tanto, valores diferentes de $\hat{\theta}_L$ y $\hat{\theta}_U$, estos puntos extremos del intervalo son valores de las variables aleatorias correspondientes $\hat{\Theta}_L$ y $\hat{\Theta}_U$. De la distribución muestral de $\hat{\Theta}$ seremos capaces de determinar $\hat{\Theta}_L$ y $\hat{\Theta}_U$, de manera que $P(\hat{\Theta}_L < \theta < \hat{\Theta}_U)$ sea igual a cualquier

valor positivo de una fracción que queremos especificar. Si, por ejemplo, calculamos $\hat{\Theta}_L$ y $\hat{\Theta}_U$ tales que

$$P(\hat{\Theta}_L < \theta < \hat{\Theta}_U) = 1 - \alpha,$$

para $0 < \alpha < 1$, tenemos entonces una probabilidad de $1 - \alpha$ de seleccionar una muestra aleatoria que produzca un intervalo que contenga θ . El intervalo $\hat{\theta}_L < \theta < \hat{\theta}_U$, que se calcula a partir de la muestra seleccionada, se llama entonces **intervalo de confianza** del $100(1 - \alpha)\%$, la fracción $1 - \alpha$ se denomina **coeficiente de confianza** o **grado de confianza**, y los extremos, $\hat{\theta}_L$ y $\hat{\theta}_U$, se denominan **límites de confianza** inferior y superior. Así, cuando $\alpha = 0.05$, tenemos un intervalo de confianza del 95%, y cuando $\alpha = 0.01$ obtenemos un intervalo de confianza más amplio del 99%. Cuanto más amplio sea el intervalo de confianza, más confiaremos en que contiene el parámetro desconocido. Desde luego, es mejor tener un 95% de confianza en que la vida promedio de cierto transistor de un televisor está entre los 6 y los 7 años, que tener un 99% de confianza en que esté entre los 3 y los 10 años. De manera ideal, preferimos un intervalo corto con un grado de confianza alto. Algunas veces las restricciones en el tamaño de nuestra muestra nos impiden tener intervalos cortos sin sacrificar cierto grado de confianza.

En las siguientes secciones estudiaremos los conceptos de estimación puntual y por intervalos, y en cada sección presentaremos un caso especial diferente. El lector debería notar que, aunque la estimación puntual y por intervalos representan diferentes aproximaciones para obtener información respecto a un parámetro, están relacionadas debido a que los estimadores del intervalo de confianza se basan en estimadores puntuales. En la siguiente sección, por ejemplo, veremos que \bar{X} es un estimador puntual de μ muy razonable. Como resultado, el importante estimador del intervalo de confianza de μ depende del conocimiento de la distribución muestral de \bar{X} .

Empezaremos la siguiente sección con el caso más sencillo de un intervalo de confianza, en donde el escenario es simple pero poco realista. Nos interesa estimar una media de la población μ cuando σ todavía se desconoce. Evidentemente, si se desconoce μ es muy improbable que se conozca σ . Cualquier información histórica que produzca datos suficientes para permitir suponer que se conoce σ probablemente habría producido información similar acerca de μ . A pesar de este argumento iniciamos con este caso porque los conceptos y los mecanismos resultantes asociados con la estimación del intervalo de confianza también estarán asociados con las situaciones más realistas que presentaremos más adelante en la sección 9.4 y las siguientes.

9.4 Una sola muestra: estimación de la media

La distribución muestral de \bar{X} está centrada en μ y en la mayoría de las aplicaciones la varianza es más pequeña que la de cualesquiera otros estimadores de μ . Por lo tanto, se utilizará la media muestral \bar{x} como una estimación puntual para la media de la población μ . Recuerde que $\sigma_{\bar{X}}^2 = \sigma^2/n$, por lo que una muestra grande producirá un valor de \bar{X} procedente de una distribución muestral con varianza pequeña. Por consiguiente, es probable que \bar{x} sea una estimación muy precisa de μ cuando n es grande.

Consideremos ahora la estimación por intervalos de μ . Si seleccionamos nuestra muestra a partir de una población normal o, a falta de ésta, si n es suficientemente grande, podemos establecer un intervalo de confianza para μ considerando la distribución muestral de \bar{X} .

De acuerdo con el teorema del límite central, podemos esperar que la distribución muestral de \bar{X} esté distribuida de forma aproximadamente normal con media $\mu_{\bar{X}} = \mu$ y desviación estándar $\sigma_{\bar{X}} = \sigma/\sqrt{n}$. Al escribir $z_{\alpha/2}$ para el valor z por arriba del cual encontramos una área de $\alpha/2$ bajo la curva normal, en la figura 9.2 podemos ver que

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha,$$

donde

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}.$$

En consecuencia,

$$P\left(-z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < z_{\alpha/2}\right) = 1 - \alpha.$$

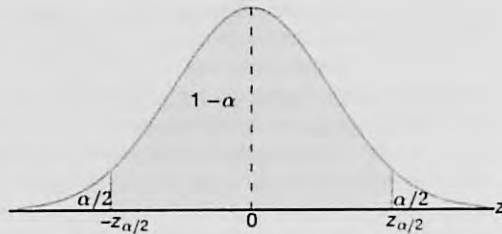


Figura 9.2: $P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha$.

Si multiplicamos cada término en la desigualdad por σ/\sqrt{n} y después restamos \bar{X} de cada término, y en seguida multiplicamos por -1 (para invertir el sentido de las desigualdades), obtenemos

$$P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

Se selecciona una muestra aleatoria de tamaño n de una población cuya varianza σ^2 se conoce y se calcula la media \bar{x} para obtener el intervalo de confianza $100(1 - \alpha)\%$. Es importante enfatizar que recurrimos al teorema del límite central citado anteriormente. Como resultado, es importante observar las condiciones para las aplicaciones que siguen.

Intervalo de confianza de μ cuando se conoce σ^2 Si \bar{x} es la media de una muestra aleatoria de tamaño n de una población de la que se conoce su varianza σ^2 , lo que da un intervalo de confianza de $100(1 - \alpha)\%$ para μ es

$$\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}},$$

donde $z_{\alpha/2}$ es el valor z que deja una área de $\alpha/2$ a la derecha.

En el caso de muestras pequeñas que se seleccionan de poblaciones no normales, no podemos esperar que nuestro grado de confianza sea preciso. Sin embargo, para muestras

de tamaño $n \geq 30$, en las que la forma de las distribuciones no esté muy sesgada, la teoría de muestreo garantiza buenos resultados.

Queda claro que los valores de las variables aleatorias $\hat{\theta}_L$ y $\hat{\theta}_U$, las cuales se definieron en la sección 9.3, son los límites de confianza

$$\hat{\theta}_L = \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \quad \text{y} \quad \hat{\theta}_U = \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

Muestras diferentes producirán valores diferentes de \bar{x} y, por lo tanto, producirán diferentes estimaciones por intervalos del parámetro μ , como se muestra en la figura 9.3. Los puntos en el centro de cada intervalo indican la posición de la estimación puntual \bar{x} para cada muestra aleatoria. Observe que todos los intervalos tienen el mismo ancho, pues esto depende sólo de la elección de $z_{\alpha/2}$ una vez que se determina \bar{x} . Cuanto más grande sea el valor de $z_{\alpha/2}$ que elijamos, más anchos haremos todos los intervalos, y podremos tener más confianza en que la muestra particular que seleccionemos producirá un intervalo que contenga el parámetro desconocido μ . En general, para una elección de $z_{\alpha/2}$, $100(1 - \alpha)\%$ de los intervalos contendrá μ .

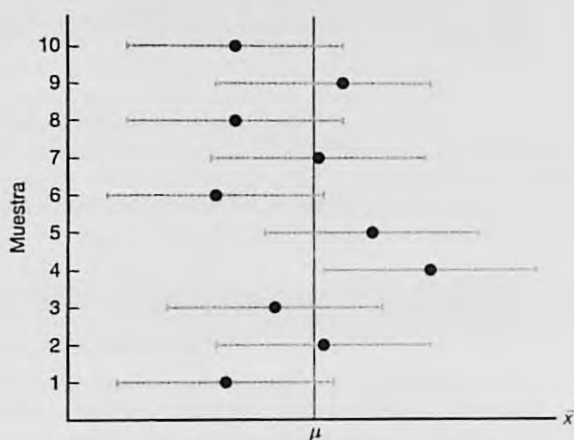


Figura 9.3: Estimaciones por intervalos de μ para muestras diferentes.

Ejemplo 9.2: Se encuentra que la concentración promedio de zinc que se obtiene en una muestra de mediciones en 36 sitios diferentes de un río es de 2.6 gramos por mililitro. Calcule los intervalos de confianza del 95% y 99% para la concentración media de zinc en el río. Suponga que la desviación estándar de la población es de 0.3 gramos por mililitro.

Solución: La estimación puntual de μ es $\bar{x} = 2.6$. El valor z que deja una área de 0.025 a la derecha y, por lo tanto, una área de 0.975 a la izquierda es $z_{0.025} = 1.96$ (véase la tabla A.3). En consecuencia, el intervalo de confianza del 95% es

$$2.6 - (1.96) \left(\frac{0.3}{\sqrt{36}} \right) < \mu < 2.6 + (1.96) \left(\frac{0.3}{\sqrt{36}} \right).$$

que se reduce a $2.50 < \mu < 2.70$. Para calcular un intervalo de confianza del 99% encontramos el valor z que deja una área de 0.005 a la derecha y de 0.995 a la izquierda. Por lo tanto, usando la tabla A.3 nuevamente, $z_{0.005} = 2.575$ y el intervalo de confianza de 99% es

$$2.6 - (2.575) \left(\frac{0.3}{\sqrt{36}} \right) < \mu < 2.6 + (2.575) \left(\frac{0.3}{\sqrt{36}} \right),$$

o simplemente

$$2.47 < \mu < 2.73.$$

Ahora vemos que se requiere un intervalo más grande para estimar μ con un mayor grado de confianza.

El intervalo de confianza del $100(1 - \alpha)\%$ ofrece un estimado de la precisión de nuestra estimación puntual. Si μ es realmente el valor central del intervalo, entonces \bar{x} estima μ sin error. La mayoría de las veces, sin embargo, \bar{x} no será exactamente igual a μ y la estimación puntual será errónea. La magnitud de este error será el valor absoluto de la diferencia entre μ y \bar{x} , de manera que podemos tener $100(1 - \alpha)\%$ de confianza en que esta diferencia no excederá a $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$. Podemos ver esto fácilmente dibujando un diagrama de un intervalo de confianza hipotético, como el de la figura 9.4.

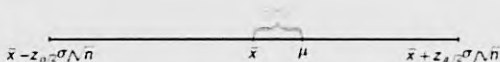


Figura 9.4: Error en la estimación de μ mediante \bar{x} .

Teorema 9.1: Si utilizamos \bar{x} como una estimación de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá a $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$.

En el ejemplo 9.2 tenemos una confianza del 95% en que la media muestral $\bar{x} = 2.6$ difiere de la media verdadera μ en una cantidad menor que $(1.96)(0.3)/\sqrt{36} = 0.1$ y 99% de confianza en que la diferencia es menor que $(2.575)(0.3)/\sqrt{36} = 0.13$.

Con frecuencia queremos saber qué tan grande necesita ser una muestra para poder estar seguros de que el error al estimar μ será menor que una cantidad específica e . Por medio del teorema 9.1 debemos elegir n de manera que $z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = e$. Al resolver esta ecuación obtenemos la siguiente fórmula para n .

Teorema 9.2: Si usamos \bar{x} como una estimación de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá a una cantidad específica e cuando el tamaño de la muestra sea

$$n = \left(\frac{z_{\alpha/2} \sigma}{e} \right)^2.$$

Cuando resolvemos para la muestra con tamaño n , redondeamos todos los valores decimales al siguiente número entero. Si seguimos este principio, podemos estar seguros de que nuestro grado de confianza nunca caerá por debajo del $100(1 - \alpha)\%$.

En términos estrictos, la fórmula del teorema 9.2 sólo será aplicable si se conoce la varianza de la población de la cual se seleccionó la muestra. Si no contamos con esa información, podríamos tomar una muestra preliminar de tamaño $n \geq 30$ para proporcionar una estimación de σ . Después, usando s como aproximación para σ en el teorema 9.2, podemos determinar aproximadamente cuántas observaciones necesitamos para brindar el grado de precisión deseado.

Ejemplo 9.3: ¿Qué tan grande debe ser la muestra del ejemplo 9.2 si queremos tener 95% de confianza en que nuestra estimación de μ diferirá por menos de 0.05?

Solución: La desviación estándar de la población es $\sigma = 0.3$. Entonces, por medio del teorema 9.2,

$$n = \left[\frac{(1.96)(0.3)}{0.05} \right]^2 = 138.3.$$

Por lo tanto, podemos tener 95% de confianza en que una muestra aleatoria de tamaño 139 proporcionará una estimación \bar{x} que diferirá de μ en una cantidad menor que 0.05. ▮

Límites de confianza unilaterales

Los intervalos de confianza y los límites de confianza resultantes que hasta ahora hemos analizado en realidad son *bilaterales*, es decir, tienen límites superior e inferior. Sin embargo, hay muchas aplicaciones en las que sólo se requiere un límite. Por ejemplo, si a un ingeniero le interesara determinar una medida de resistencia a la tensión, la información que más le ayudaría a lograr su objetivo sería la del límite inferior, ya que éste indica el escenario del “peor caso”, es decir, el de la menor resistencia. Por otro lado, si se buscara determinar una medida para la cual un valor de μ relativamente grande no fuera redituable o deseable, entonces la medida que resultaría de interés sería la del límite de confianza superior. Un ejemplo en el que la medida del límite superior sería muy informativa es el caso en el que se necesita hacer inferencias para determinar la composición media de mercurio en el agua de un río.

Los límites de confianza unilaterales se desarrollan de la misma forma que los intervalos bilaterales. Sin embargo, la fuente es un enunciado de probabilidad unilateral que utiliza el teorema del límite central:

$$P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < z_\alpha\right) = 1 - \alpha.$$

Entonces, es posible manipular el enunciado de probabilidad de forma muy similar a como se hizo anteriormente para obtener

$$P(\mu > \bar{X} - z_\alpha \sigma/\sqrt{n}) = 1 - \alpha.$$

Una manipulación similar de $P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} > -z_\alpha\right) = 1 - \alpha$ da

$$P(\mu < \bar{X} + z_\alpha \sigma/\sqrt{n}) = 1 - \alpha.$$

Como resultado, se obtienen los siguientes límites unilaterales superior e inferior.

Límites de confianza unilaterales de μ cuando se conoce el valor de σ^2	Si \bar{X} es la media de una muestra aleatoria de tamaño n a partir de una población con varianza σ^2 , los límites de confianza unilaterales del $100(1 - \alpha)\%$ para μ son dados por
	límite unilateral superior: $\bar{x} + z_\alpha \sigma/\sqrt{n}$;
	límite unilateral inferior: $\bar{x} - z_\alpha \sigma/\sqrt{n}$.

Ejemplo 9.4: En un experimento de pruebas psicológicas se seleccionan al azar 25 sujetos y se miden sus tiempos de reacción, en segundos, ante un estímulo particular. La experiencia sugiere que la varianza en los tiempos de reacción ante los diferentes tipos de estímulos es de 4 s^2 y que la distribución del tiempo de reacción es aproximadamente normal. El tiempo promedio para los sujetos fue de 6.2 segundos. Calcule un límite superior del 95% para el tiempo medio de reacción.

Solución: Lo que da el límite superior del 95% es

$$\begin{aligned}\bar{x} + z_{\alpha} \sigma / \sqrt{n} &= 6.2 + (1.645) \sqrt{4/25} = 6.2 + 0.658 \\ &= 6.858 \text{ segundos.}\end{aligned}$$

En consecuencia, tenemos un 95% de confianza en que el tiempo promedio de reacción es menor que 6.858 segundos. J

El caso en que se desconoce σ

Con frecuencia debemos tratar de estimar la media de una población sin conocer la varianza. El lector debería recordar que en el capítulo 8 aprendió que, si tenemos una muestra aleatoria a partir de una *distribución normal*, entonces la variable aleatoria

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

tiene una distribución *t* de Student con $n - 1$ grados de libertad. Aquí S es la desviación estándar de la muestra. En esta situación, en la que se desconoce σ , se puede utilizar T para construir un intervalo de confianza para μ . El procedimiento es igual que cuando se conoce σ , sólo que en este caso σ se reemplaza con S y la distribución normal estándar se reemplaza con la distribución *t*. Si nos remitimos a la figura 9.5, podemos afirmar que

$$P(-t_{\alpha/2} < T < t_{\alpha/2}) = 1 - \alpha,$$

donde $t_{\alpha/2}$ es el valor *t* con $n - 1$ grados de libertad, por arriba del cual encontramos una área de $\alpha/2$. Debido a la simetría, un área igual de $\alpha/2$ caerá a la izquierda de $-t_{\alpha/2}$. Al sustituir por T escribimos

$$P\left(-t_{\alpha/2} < \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{\alpha/2}\right) = 1 - \alpha.$$

Al multiplicar cada término en la desigualdad por S/\sqrt{n} y después restar \bar{X} de cada término y multiplicar por -1 , obtenemos

$$P\left(\bar{X} - t_{\alpha/2} \frac{S}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{S}{\sqrt{n}}\right) = 1 - \alpha.$$

Para nuestra muestra aleatoria particular de tamaño n se calculan la media \bar{x} y la desviación estándar s , y se obtiene el siguiente intervalo de confianza $100(1 - \alpha)\%$ para μ

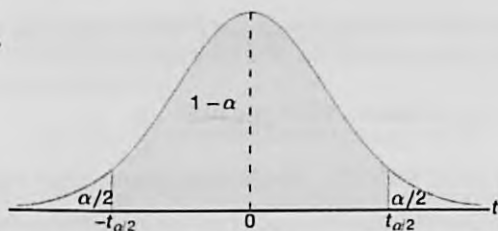


Figura 9.5: $P(-t_{\alpha/2} < T < t_{\alpha/2}) = 1 - \alpha$.

Intervalo de confianza para μ cuando se desconoce σ^2 Si \bar{x} y s son la media y la desviación estándar de una muestra aleatoria de una población normal de la que se desconoce la varianza σ^2 , un intervalo de confianza del $100(1 - \alpha)\%$ para μ es

$$\bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}},$$

donde $t_{\alpha/2}$ es el valor t con $v = n - 1$ grados de libertad que deja una área de $\alpha/2$ a la derecha.

Hicimos una distinción entre los casos en los que se conoce σ y en los que se desconoce calculando las estimaciones del intervalo de confianza. Deberíamos resaltar que para el caso en que se conoce σ se utiliza el teorema del límite central, mientras que, para el caso en que se desconoce, se usa la distribución muestral de la variable aleatoria T . Sin embargo, el uso de la distribución t se basa en la premisa de que el muestreo es de una distribución normal. Siempre que la forma de la distribución se aproxime a la de campana, se puede utilizar la distribución t para calcular los intervalos de confianza cuando se desconoce σ^2 , y se pueden esperar muy buenos resultados.

Los límites de confianza unilaterales calculados para μ con σ desconocida son como el lector esperaría, a saber:

$$\bar{x} + t_{\alpha} \frac{s}{\sqrt{n}} \quad \text{y} \quad \bar{x} - t_{\alpha} \frac{s}{\sqrt{n}}.$$

Éstos son, respectivamente, los límites superior e inferior del $100(1 - \alpha)\%$. Aquí t_{α} es el valor t que tiene una área α a la derecha.

Ejemplo 9.5: El contenido de ácido sulfúrico de 7 contenedores similares es de 9.8, 10.2, 10.4, 9.8, 10.0, 10.2, y 9.6 litros. Calcule un intervalo de confianza del 95% para el contenido promedio de todos los contenedores suponiendo una distribución aproximadamente normal.

Solución: La media muestral y la desviación estándar para los datos dados son

$$\bar{x} = 10.0 \quad \text{y} \quad s = 0.283.$$

Si usamos la tabla A.4, encontramos $t_{0.025} = 2.447$ para $v = 6$ grados de libertad. En consecuencia, el intervalo de confianza del 95% para μ es

$$10.0 - (2.447) \left(\frac{0.283}{\sqrt{7}} \right) < \mu < 10.0 + (2.447) \left(\frac{0.283}{\sqrt{7}} \right),$$

que se reduce a $9.74 < \mu < 10.26$.

Concepto de intervalo de confianza para una muestra grande

Con frecuencia los estadísticos recomiendan que incluso cuando no sea posible suponer la normalidad, se desconozca σ y $n \geq 30$, σ se puede reemplazar con s para poder utilizar el intervalo de confianza

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$$

A menudo se hace referencia a esto como un *intervalo de confianza para una muestra grande*. La justificación para esto reside sólo en la presunción de que, con una muestra tan grande como 30 y una distribución de la población no muy sesgada, s estará muy cerca de la σ verdadera y, de esta manera, el teorema del límite central continuará siendo válido. Se debería destacar que esto es sólo una aproximación y que la calidad de los resultados mejora a medida que aumenta el tamaño de la muestra.

Ejemplo 9.6: Se obtienen las calificaciones de matemáticas del Examen de Aptitudes Escolares (SAT, por sus siglas en inglés) de una muestra aleatoria de 500 estudiantes del último año de preparatoria del estado de Texas. Se calculan la media y la desviación estándar muestrales, que son 501 y 112, respectivamente. Calcule un intervalo de confianza del 99% de la calificación promedio de matemáticas en el SAT para los estudiantes del último año de preparatoria del estado de Texas.

Solución: Como el tamaño de la muestra es grande, es razonable utilizar la aproximación normal. Si utilizamos la tabla A.3, encontramos $z_{0.005} = 2.575$. Por lo tanto, un intervalo de confianza del 99% para μ es

$$501 \pm (2.575) \left(\frac{112}{\sqrt{500}} \right) = 501 \pm 12.9,$$

que da como resultado $488.1 < \mu < 513.9$.

9.5 Error estándar de una estimación puntual

Hicimos una distinción muy clara entre los objetivos de las estimaciones puntuales y las estimaciones del intervalo de confianza. Las primeras proporcionan un solo número que se extrae de un conjunto de datos experimentales, y las segundas proporcionan un intervalo razonable para el parámetro, *dados los datos experimentales*; es decir, $100(1 - \alpha)\%$ de tales intervalos que se calcula “cubren” el parámetro.

Estos dos métodos de estimación se relacionan entre sí. El elemento en común es la distribución muestral del estimador puntual. Considere, por ejemplo, el estimador \bar{X} de μ cuando se conoce σ . Indicamos antes que una medida de la calidad de un estimador insesgado es su varianza. La varianza de \bar{X} es

$$\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}.$$

De esta forma, la desviación estándar de \bar{X} o *error estándar de \bar{X}* es σ/\sqrt{n} . En términos simples, el error estándar de un estimador es su desviación estándar. Para el caso de \bar{X} el límite de confianza que se calcula

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \text{ se escribe como } \bar{x} \pm z_{\alpha/2} \text{ e.e. } (\bar{x}),$$

donde "e.e." es el error estándar. El punto importante es que el ancho del intervalo de confianza de μ depende de la calidad del estimador puntual a través de su error estándar. En el caso en que se desconoce σ y la muestra proviene de una distribución normal, s reemplaza a σ y se incluye el *error estándar estimado* S/\sqrt{n} . Por consiguiente, los límites de confianza de μ son:

Límites de
confianza para μ
cuando se
desconoce σ^2

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} = \bar{x} \pm t_{\alpha/2} \text{ e.e.}(\bar{x})$$

De nuevo, el intervalo de confianza *no es mejor* (en términos de anchura) *que la calidad de la estimación puntual*, en este caso a través de su error estándar estimado. A menudo el software de computación se refiere a los errores estándar estimados simplemente como "errores estándar".

A medida que avanzamos a intervalos de confianza más complejos, prevalece el concepto de que el ancho de los intervalos de confianza se acorta cuando mejora la calidad de la estimación puntual correspondiente, aunque no siempre es tan sencillo como aquí se ilustra. Se puede argumentar que un intervalo de confianza es tan sólo una ampliación de la estimación puntual para tomar en cuenta la exactitud de dicha estimación.

9.6 Intervalos de predicción

La estimación puntual y la estimación por intervalos de la media que se expusieron en las secciones 9.4 y 9.5 proporcionan buena información del parámetro desconocido μ de una distribución normal, o de una distribución no normal a partir de la cual se toma una muestra grande. Algunas veces, además de la media de la población, el experimentador podría estar interesado en predecir el **valor posible de una observación futura**. Por ejemplo, en el control de calidad el experimentador podría necesitar utilizar los datos observados para predecir una nueva observación. Un proceso de manufactura de una pieza de metal se podría evaluar basándose en si la pieza cumple con las especificaciones de resistencia a la tensión. En ciertas ocasiones un cliente podría estar interesado en comprar una **sola pieza**. En este caso un intervalo de confianza de la resistencia media a la tensión no cubriría la información requerida. El cliente necesitaría una aseveración respecto a la incertidumbre de una **sola observación**. Este tipo de requerimiento se satisface muy bien construyendo un **intervalo de predicción**.

Es muy sencillo obtener un intervalo de predicción para las situaciones que hemos considerado hasta el momento. Suponga que la muestra aleatoria se tomó de una población normal con media μ desconocida y varianza σ^2 conocida. Un estimador puntual natural de una nueva observación es \bar{X} . En la sección 8.4 se aprendió que la varianza de \bar{X} es σ^2/n . Sin embargo, para predecir una nueva observación no basta con explicar la variación debida a la estimación de la media, también tendríamos que explicar la **variación de una observación futura**. A partir de la suposición sabemos que la varianza del

error aleatorio en una nueva observación es σ^2 . El desarrollo de un intervalo de predicción se representa mejor empezando con una variable aleatoria normal $x_0 - \bar{x}$, donde x_0 es la nueva observación y \bar{x} se toma de la muestra. Como x_0 y \bar{x} son independientes, sabemos que

$$z = \frac{x_0 - \bar{x}}{\sqrt{\sigma^2 + \sigma^2/n}} = \frac{x_0 - \bar{x}}{\sigma\sqrt{1 + 1/n}}$$

es $n(z; 0, 1)$. Como resultado, si utilizamos el enunciado de probabilidad

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha$$

con el estadístico z anterior, y si colocamos x_0 en el centro del enunciado de probabilidad, tenemos que la probabilidad de que ocurra el siguiente evento es $1 - \alpha$:

$$\bar{x} - z_{\alpha/2}\sigma\sqrt{1 + 1/n} < x_0 < \bar{x} + z_{\alpha/2}\sigma\sqrt{1 + 1/n}$$

Como resultado, el intervalo de predicción calculado se formaliza como sigue.

Intervalo de predicción para una observación futura cuando se conoce σ^2

Para una distribución normal de mediciones con media μ desconocida y varianza σ^2 conocida, un **intervalo de predicción** del $100(1 - \alpha)\%$ de una observación futura x_0 es

$$\bar{x} - z_{\alpha/2}\sigma\sqrt{1 + 1/n} < x_0 < \bar{x} + z_{\alpha/2}\sigma\sqrt{1 + 1/n},$$

donde $z_{\alpha/2}$ es el valor z que deja una área de $\alpha/2$ a la derecha.

Ejemplo 9.7: Debido a la disminución en las tasas de interés el First Citizens Bank recibió muchas solicitudes para hipoteca. Una muestra reciente de 50 créditos hipotecarios dio como resultado un promedio en la cantidad de préstamos de \$257,300. Suponga una desviación estándar de la población de \$25,000. En el caso del siguiente cliente que llena una solicitud de crédito hipotecario calcule un intervalo de predicción del 95% para la cantidad del crédito.

Solución: La predicción puntual de la cantidad del crédito del siguiente cliente es $\bar{x} = \$257,300$. El valor z aquí es $z_{0.025} = 1.96$. Por lo tanto, un intervalo de predicción del 95% para la cantidad de un crédito futuro es

$$257,300 - (1.96)(25,000)\sqrt{1 + 1/50} < x_0 < 257,300 + (1.96)(25,000)\sqrt{1 + 1/50},$$

que produce el intervalo (\$207,812.43, \$306,787.57). ▮

El intervalo de predicción proporciona un buen estimado de la ubicación de una observación futura, el cual es muy diferente del estimado del valor promedio de la muestra. Debe advertirse que la variación de esta predicción es la suma de la variación debida a una estimación de la media y la variación de una sola observación. Sin embargo, como antes, consideramos primero el caso en el que se conoce la varianza. En el caso en que se desconoce la varianza también es importante tratar con el intervalo de predicción de una observación futura. De hecho, en este caso se podría utilizar una distribución t de Student, como se describe en el siguiente resultado. Aquí la distribución normal simplemente se reemplaza con la distribución t .

Intervalo de predicción de una observación futura cuando se desconoce σ^2

Para una distribución normal de mediciones cuando la media μ y la varianza σ^2 se desconocen, un **intervalo de predicción** del $100(1 - \alpha)\%$ de una observación futura x_0 es

$$\bar{x} - t_{\alpha/2} s \sqrt{1 + 1/n} < x_0 < \bar{x} + t_{\alpha/2} s \sqrt{1 + 1/n},$$

donde $t_{\alpha/2}$ es el valor t con $v = n - 1$ grados de libertad, que deja una área de $\alpha/2$ a la derecha.

También se pueden utilizar intervalos de predicción unilaterales. Los límites de predicción superiores se aplican en casos en los que es necesario enfocarse en observaciones futuras grandes. El interés por observaciones pequeñas futuras requiere utilizar límites de predicción más bajos. El límite superior es dado por

$$\bar{x} + t_{\alpha} s \sqrt{1 + 1/n}$$

y el límite inferior por

$$\bar{x} - t_{\alpha} s \sqrt{1 + 1/n}.$$

Ejemplo 9.8: Un inspector de alimentos seleccionó aleatoriamente 30 paquetes de carne de res 95% magra. La muestra dio como resultado una media de 96.2% con una desviación estándar muestral de 0.8%. Calcule un intervalo de predicción del 99% para la condición baja en grasa de un paquete nuevo. Suponga normalidad.

Solución: Para $v = 29$ grados de libertad, $t_{0.005} = 2.756$. Por lo tanto, un intervalo de predicción del 99% para una observación nueva x_0 es

$$96.2 - (2.756)(0.8) \sqrt{1 + \frac{1}{30}} < x_0 < 96.2 + (2.756)(0.8) \sqrt{1 + \frac{1}{30}},$$

que se reduce a (93.96, 98.44). ▮

Uso de límites de predicción para detectar valores extremos

Hasta el momento hemos puesto poca atención al concepto de **valores extremos** u observaciones aberrantes. La mayoría de los investigadores científicos son muy sensibles a la existencia de observaciones de valores extremos, también llamados datos defectuosos o “malos”. En el capítulo 12 profundizaremos en el estudio de este concepto. Sin embargo, nos interesa considerarlos aquí porque la detección de los valores extremos está estrechamente relacionada con los intervalos de predicción.

Para nuestros propósitos nos conviene considerar que una observación extrema es una que proviene de una población con una media diferente a la que determina el resto de la muestra de tamaño n que se está estudiando. El intervalo de predicción produce un límite que “cubre” una sola observación futura con probabilidad $1 - \alpha$, si ésta proviene de la población de la que se tomó la muestra. Por lo tanto, una metodología para detectar valores extremos implica la regla de que **una observación es un valor extremo si cae fuera del intervalo de predicción calculado sin incluir la observación cuestionable en la muestra**. Como resultado, para el intervalo de predicción del ejemplo 9.8, en el caso de los paquetes de carne, la observación que se obtiene al medir un nuevo paquete y encontrar que su contenido libre de grasa está fuera del intervalo (93.96, 98.44) se podría considerar como un valor extremo.

9.7 Límites de tolerancia

Como vimos en la sección 9.6, el científico o el ingeniero podrían estar menos interesados en estimar parámetros que en obtener información sobre el lugar en el que caería una *observación* o medición individual. Este tipo de situaciones requiere intervalos de predicción. Sin embargo, existe un tercer tipo de intervalo que es útil en muchas aplicaciones. Una vez más, suponga que el interés se centra en torno a la fabricación de la pieza de un componente y que existen especificaciones sobre una dimensión de esa parte. Además, la media de esa dimensión no es tan importante. Sin embargo, a diferencia del escenario de la sección 9.6, se podría estar menos interesado en una sola observación y más en el lugar en el que cae la mayoría de la población. Si las especificaciones del proceso son importantes, el administrador del proceso se interesará en el desempeño a largo plazo, **no en la siguiente observación**. Debemos tratar de determinar los límites que, en cierto sentido probabilístico, "cubren" los valores en la población, es decir, los valores medidos de la dimensión.

Un método para establecer el límite deseado consiste en determinar un intervalo de confianza sobre una *proporción fija* de las mediciones. Esto se comprende mejor visualizando una situación en la que se realiza un muestreo aleatorio de una distribución normal con media conocida μ y varianza σ^2 . Evidentemente, un límite que cubre el 95% central de la población de observaciones es

$$\mu \pm 1.96\sigma.$$

A esto se le llama **intervalo de tolerancia** y, en realidad, su cobertura del 95% de las observaciones medidas es exacta. Sin embargo, en la práctica rara vez se conocen μ y σ ; por consiguiente, el usuario debe aplicar

$$\bar{x} \pm ks.$$

Ahora bien, el intervalo es, desde luego, una variable aleatoria, por lo tanto, la *cobertura* de una proporción de la población por el intervalo no es exacta. Como resultado, se debe usar un intervalo de confianza del $100(1 - \gamma)\%$, ya que no se puede esperar que $\bar{x} \pm ks$ cubra cualquier proporción específica todo el tiempo. Lo anterior nos lleva a la siguiente definición.

Límites de tolerancia Para una distribución normal de mediciones en la que se desconoce la media μ y la desviación estándar σ , los **límites de tolerancia** son dados por $\bar{x} \pm ks$, donde k se determina de tal manera que se pueda estar seguro, con un $100(1 - \gamma)\%$ de confianza, de que los límites dados contienen al menos la proporción $1 - \alpha$ de las mediciones.

La tabla A.7 ofrece valores de k para $1 - \alpha = 0.90, 0.95, 0.99$; $\gamma = 0.05, 0.01$; y para valores seleccionados de n de 2 a 300.

Ejemplo 9.9: Considere el ejemplo 9.8. Con la información dada calcule un intervalo de tolerancia que proporcione límites bilaterales del 95% sobre el 90% de la distribución de paquetes de carne 95% magra. Suponga que los datos provienen de una distribución aproximadamente normal.

Solución: Del ejemplo 9.8, recuerde que $n = 30$, que la media muestral es de 96.2% y que la desviación estándar muestral es de 0.8%. De la tabla A.7, $k = 2.14$. Si utilizamos

$$\bar{x} \pm ks = 96.2 \pm (2.14)(0.8),$$

encontramos que los límites inferior y superior son de 94.5 y de 97.9.

Tenemos 95% de confianza en que el rango anterior cubre el 90% central de la distribución de paquetes de carne de res 95% magra.

Diferencia entre intervalos de confianza, intervalos de predicción e intervalos de tolerancia

Es importante resaltar la diferencia entre los tres tipos de intervalos que se estudiaron e ilustraron en las secciones anteriores. Los cálculos son sencillos, pero la interpretación podría resultar confusa. En aplicaciones de la vida real tales intervalos no son intercambiables, ya que sus interpretaciones son muy diferentes.

En el caso de los intervalos de confianza sólo se pone atención en la **media de la población**. Por ejemplo, el ejercicio 9.13 de la página 283 se refiere a un proceso de ingeniería que produce alfileres para costura. Se establece una especificación sobre la dureza de Rockwell por debajo de la cual el cliente no aceptará ningún alfiler. En este caso un parámetro de la población debe tener poca relevancia. Es importante que el ingeniero sepa en dónde *van a estar la mayoría de los valores de la dureza de Rockwell*. Por consiguiente, se deberían utilizar los límites de tolerancia. Seguramente, al administrador le agrada saber que los límites de tolerancia en cualquier producto del proceso son más rigurosos que las especificaciones para el propio proceso.

Es verdad que la interpretación del límite de tolerancia se relaciona hasta cierto punto con el intervalo de confianza. El intervalo de tolerancia del $100(1 - \alpha)\%$ sobre, digamos, la proporción 0.95, se podría considerar como un intervalo de confianza **sobre el 95% intermedio** de la distribución normal correspondiente. Los límites de tolerancia unilaterales también son relevantes. En el caso del problema de dureza de Rockwell se desearía tener un límite inferior de la forma $\bar{x} - ks$, tal que se tenga un 99% de confianza en que al menos 99% de los valores de la dureza de Rockwell excederán al valor calculado.

Los intervalos de predicción se pueden aplicar cuando es importante determinar un límite para un **solo valor**. Aquí la media no es la cuestión, ni tampoco la ubicación de la mayoría de la población, lo que se requiere, más bien, es la ubicación de una sola nueva observación.

Estudio de caso 9.1: Calidad de una máquina. Una máquina produce piezas de metal que tienen forma cilíndrica. Se toma una muestra de tales piezas y se encuentra que los diámetros son 1.01, 0.97, 1.03, 1.04, 0.99, 0.98, 0.99, 1.01 y 1.03 centímetros. Utilice estos datos para calcular tres tipos de intervalos y hacer interpretaciones que ilustren las diferencias entre ellos en el contexto del sistema. Para todos los cálculos suponga una distribución aproximadamente normal. La media muestral y la desviación estándar para los datos dados son $\bar{x} = 1.0056$ y $s = 0.0246$.

- Calcule un intervalo de confianza del 99% sobre la media del diámetro.
- Calcule un intervalo de predicción del 99% sobre el diámetro medido de una sola pieza de metal tomada de la máquina.
- Calcule los límites de tolerancia del 99% que contengan 95% de las piezas de metal producidas por esta máquina.

Solución: a) El intervalo de confianza del 99% para la media del diámetro está dado por

$$\bar{x} \pm t_{0.005, s/\sqrt{n}} = 1.0056 \pm (3.355)(0.0246/3) = 1.0056 \pm 0.0275.$$

Por lo tanto, los límites de confianza del 99% son 0.9781 y 1.0331.

- b) El intervalo de predicción del 99% para una futura observación está dado por

$$\bar{x} \pm t_{0.005} s \sqrt{1 + 1/n} = 1.0056 \pm (3.355)(0.0246) \sqrt{1 + 1/9},$$

donde los límites son 0.9186 y 1.0926.

- c) De la tabla A.7, para $n = 9$, $1 - \gamma = 0.99$, y $1 - \alpha = 0.95$, obtenemos $k = 4.550$ para los límites bilaterales. Por lo tanto, los límites de tolerancia del 99% son dados por

$$\bar{x} + ks = 1.0056 \pm (4.550)(0.0246),$$

donde los límites son 0.8937 y 1.1175. Tenemos un 99% de confianza en que el intervalo de tolerancia de 0.8937 a 1.1175 contendrá el 95% central de la distribución de diámetros producidos.

Este estudio de caso ilustra que los tres tipos de límites pueden conducir a resultados muy diferentes, aunque todos son límites del 99%. En el caso del intervalo de confianza sobre la media, el 99% de estos intervalos cubre la media del diámetro de la población. Por lo tanto, decimos que tenemos un 99% de confianza en que la media del diámetro producido por el proceso se encuentra entre 0.9781 y 1.0331 centímetros. Se hace hincapié en la media y se pone poco interés en una sola lectura o en la naturaleza general de la distribución de diámetros en la población. En lo que se refiere a los límites de predicción, los límites 0.9186 y 1.0926 se basan en la distribución de una sola pieza "nueva" de metal tomada del proceso, y nuevamente el 99% de estos límites cubren el diámetro de una nueva pieza medida. Por otro lado, como se sugirió en la sección anterior, los límites de tolerancia le dan al ingeniero una idea de en qué parte de la población se localiza la "mayoría", digamos el 95% central, de los diámetros de las piezas medidas. Los límites de tolerancia del 99%, 0.8937 y 1.1175 difieren mucho de los otros dos límites. Si esos límites le parecen demasiado anchos al ingeniero, esto se reflejará de forma negativa en la calidad del proceso. Por otro lado, si los límites representan un resultado deseable, el ingeniero podría concluir que la mayoría (95% en este caso) de los diámetros se encuentran dentro de un rango adecuado. De nuevo, se podría hacer una interpretación del intervalo de confianza, a saber, el 99% de esos límites calculados cubrirán el 95% intermedio de la población de diámetros. J

Ejercicios

9.1 Un investigador de la UCLA afirma que la esperanza de vida de los ratones se puede extender hasta en 25% cuando se reduce aproximadamente 40% de las calorías de su dieta desde el momento en que son destetados. La dieta restringida se enriquece hasta niveles normales con vitaminas y proteínas. Si se supone que a partir de estudios previos se sabe que $\sigma = 5.8$ meses, ¿cuántos ratones se deberían incluir en la muestra para tener un 99% de confianza en que la vida media esperada de la muestra estará dentro de 2 meses a partir de la media de la población para todos los ratones sujetos a la dieta reducida?

9.2 Una empresa de material eléctrico fabrica bombillas que tienen una duración distribuida de forma aproximadamente normal, con una desviación estándar

de 40 horas. Si una muestra de 30 bombillas tiene una duración promedio de 780 horas, calcule un intervalo de confianza del 96% para la media de la población de todas las bombillas producidas por esta empresa.

9.3 Muchos pacientes con problemas del corazón tienen un marcapasos para controlar su ritmo cardiaco. El marcapasos tiene montado un módulo conector de plástico en la parte superior. Suponga una desviación estándar de 0.0015 pulgadas y una distribución aproximadamente normal, y con base en esto calcule un intervalo de confianza del 95% para la media de la profundidad de todos los módulos conectores fabricados por cierta empresa. Una muestra aleatoria de 75 módulos tiene una profundidad promedio de 0.310 pulgadas.

9.4 Las estaturas de una muestra aleatoria de 50 estudiantes universitarios tienen una media de 174.5 centímetros y una desviación estándar de 6.9 centímetros.

- Construya un intervalo de confianza del 98% para la estatura media de todos los estudiantes universitarios.
- ¿Qué podemos afirmar con una confianza del 98% acerca del posible tamaño de nuestro error, si estimamos que la estatura media de todos los estudiantes universitarios es de 174.5 centímetros?

9.5 Una muestra aleatoria de 100 propietarios de automóviles del estado de Virginia revela que éstos conducen su automóvil, en promedio, 23,500 kilómetros por año, con una desviación estándar de 3900 kilómetros. Suponga que la distribución de las mediciones es aproximadamente normal.

- Construya un intervalo de confianza del 99% para el número promedio de kilómetros que un propietario de un automóvil conduce anualmente en Virginia.
- ¿Qué podemos afirmar con un 99% de confianza acerca del posible tamaño del error, si estimamos que los propietarios de automóviles de Virginia conducen un promedio de 23,500 kilómetros por año?

9.6 ¿Qué tan grande debe ser la muestra en el ejercicio 9.2 si deseamos tener un 96% de confianza en que nuestra media muestral estará dentro de 10 horas a partir de la media verdadera?

9.7 ¿De qué tamaño debe ser la muestra en el ejercicio 9.3 si deseamos tener un 95% de confianza en que nuestra media muestral estará dentro de un 0.0005 de pulgada de la media verdadera?

9.8 Un experto en eficiencia desea determinar el tiempo promedio que toma perforar tres hoyos en cierta placa metálica. ¿De qué tamaño debe ser una muestra para tener un 95% de confianza en que esta media muestral estará dentro de 15 segundos de la media verdadera? Suponga que por estudios previos se sabe que $\sigma = 40$ segundos.

9.9 Según estudios realizados por el doctor W. H. Bowen, del Instituto Nacional de Salud, y por el doctor J. Yudben, profesor de nutrición y dietética de la Universidad de Londres, el consumo regular de cereales preendulzados contribuye al deterioro de los dientes, a las enfermedades cardíacas y a otras enfermedades degenerativas. En una muestra aleatoria de 20 porciones sencillas similares del cereal Alpha-Bits, el contenido promedio de azúcar era de 11.3 gramos con una desviación estándar de 2.45 gramos. Suponga que el contenido de azúcar está distribuido normalmente y con base en esto construya un intervalo de confianza de 95% para el contenido medio de azúcar de porciones sencillas de Alpha-Bits.

9.10 Las integrantes de una muestra aleatoria de 12 graduadas de cierta escuela para secretarías teclearon

un promedio de 79.3 palabras por minuto, con una desviación estándar de 7.8 palabras por minuto. Suponga una distribución normal para el número de palabras que teclean por minuto y con base en esto calcule un intervalo de confianza del 95% para el número promedio de palabras que teclean todas las graduadas de esta escuela.

9.11 Una máquina produce piezas metálicas de forma cilíndrica. Se toma una muestra de las piezas y los diámetros son 1.01, 0.97, 1.03, 1.04, 0.99, 0.98, 0.99, 1.01 y 1.03 centímetros. Calcule un intervalo de confianza del 99% para la media del diámetro de las piezas que se manufacturan con esta máquina. Suponga una distribución aproximadamente normal.

9.12 Una muestra aleatoria de 10 barras energéticas de chocolate de cierta marca tiene, en promedio, 230 calorías por barra y una desviación estándar de 15 calorías. Construya un intervalo de confianza del 99% para el contenido medio verdadero de calorías de esta marca de barras energéticas de chocolate. Suponga que la distribución del contenido calórico es aproximadamente normal.

9.13 En un estudio para determinar la dureza de Rockwell en la cabeza de alfileres para costura se toma una muestra aleatoria de 12. Se toman mediciones de la dureza de Rockwell para cada una de las 12 cabezas y se obtiene un valor promedio de 48.50, con una desviación estándar muestral de 1.5. Suponga que las mediciones se distribuyen de forma normal y con base en esto construya un intervalo de confianza de 90% para la dureza media de Rockwell.

9.14 Se registran las siguientes mediciones del tiempo de secado, en horas, de cierta marca de pintura vinílica:

3.4	2.5	4.8	2.9	3.6
2.8	3.3	5.6	3.7	2.8
4.4	4.0	5.2	3.0	4.8

Suponga que las mediciones representan una muestra aleatoria de una población normal y con base en esto calcule el intervalo de predicción del 95% para el tiempo de secado de la siguiente prueba de pintura.

9.15 Remítase al ejercicio 9.5 y construya un intervalo de predicción del 99% para los kilómetros que viaja anualmente el propietario de un automóvil en Virginia.

9.16 Considere el ejercicio 9.10 y calcule el intervalo de predicción del 95% para el siguiente número observado de palabras por minuto tecleadas por una graduada de la escuela de secretarías.

9.17 Considere el ejercicio 9.9 y calcule un intervalo de predicción del 95% para el contenido de azúcar de la siguiente porción de cereal Alpha-Bits.

9.18 Remítase al ejercicio 9.13 y construya un intervalo de tolerancia del 95% que contenga el 90% de las mediciones.

9.19 Una muestra aleatoria de 25 tabletas de aspirina con antiácido contiene, en promedio, 325.05 mg de aspirina en cada tableta, con una desviación estándar de 0.5 mg. Calcule los límites de tolerancia del 95% que contendrán 90% del contenido de aspirina para esta marca. Suponga que el contenido de aspirina se distribuye normalmente.

9.20 Considere la situación del ejercicio 9.11. Aunque la estimación de la media del diámetro es importante, no es ni con mucho tan importante como intentar determinar la ubicación de la mayoría de la distribución de los diámetros. Calcule los límites de tolerancia del 95% que contengan el 95% de los diámetros.

9.21 En un estudio realizado por el Departamento de Zoología del Virginia Tech con el fin de conocer la cantidad de ortofósforo en el río, se recolectaron 15 "muestras" de agua en una determinada estación ubicada en el río James. La concentración del químico se midió en miligramos por litro. Suponga que la media en la estación de muestreo no es tan importante como la distribución de las concentraciones del químico en los extremos superiores. El interés se centra en saber si las concentraciones en estos extremos son demasiado elevadas. Las lecturas de las 15 muestras de agua proporcionaron una media muestral de 3.84 miligramos por litro y una desviación estándar muestral de 3.07 miligramos por litro. Suponga que las lecturas son una muestra aleatoria de una distribución normal. Calcule un intervalo de predicción (límite de predicción superior del 95%) y un límite de tolerancia (un límite de tolerancia superior del 95% que excede al 95% de la población de valores). Interprete ambos límites, es decir, especifique qué indica cada uno acerca de los extremos superiores de la distribución de ortofósforo en la estación de muestreo.

9.22 Se están estudiando las propiedades de resistencia a la tensión de un determinado tipo de hilo. Con ese fin se prueban 50 piezas en condiciones similares y los resultados que se obtienen revelan una resistencia a la tensión promedio de 78.3 kilogramos y una desviación estándar de 5.6 kilogramos. Suponga que la resistencia a la tensión tiene una distribución normal y con base en esto calcule un límite de predicción inferior al 95% de un solo valor observado de resistencia a la tensión. Además, determine un límite inferior de tolerancia del 95% que sea excedido por el 99% de los valores de resistencia a la tensión.

9.23 Remítase al ejercicio 9.22. ¿Por qué las 1/2 cantidades solicitadas en el ejercicio parecen ser más importantes para el fabricante del hilo que, por ejemplo, un intervalo de confianza en la resistencia media a la tensión?

9.24 Remítase una vez más al ejercicio 9.22. Suponga que un comprador del hilo específica que éste debe

tener una resistencia a la tensión de por lo menos 62 kilogramos. El fabricante estará satisfecho si la cantidad de piezas producidas que no cumplen la especificación no excede al 5%. ¿Hay alguna razón para preocuparse? Esta vez utilice un límite de tolerancia unilateral del 99% que sea excedido por el 95% de los valores de resistencia a la tensión.

9.25 Considere las mediciones del tiempo de secado del ejercicio 9.14. Suponga que las 15 observaciones en el conjunto de datos también incluyen un decimosexto valor de 6.9 horas. En el contexto de las 15 observaciones originales, ¿el valor decimosexto es un valor extremo? Muestre el procedimiento.

9.26 Considere los datos del ejercicio 9.13. Suponga que el fabricante de los alfileres insiste en que la dureza de Rockwell del producto es menor o igual que 44.0 sólo un 5% de las veces. ¿Cuál es su reacción? Utilice un cálculo de un límite de tolerancia como la base de su veredicto.

9.27 Considere la situación del estudio de caso 9.1 de la página 281, con una muestra más grande de piezas metálicas. Los diámetros son los siguientes: 1.01, 0.97, 1.03, 1.04, 0.99, 0.98, 1.01, 1.03, 0.99, 1.00, 1.00, 0.99, 0.98, 1.01, 1.02, 0.99 centímetros. Nuevamente puede suponer una distribución normal. Haga lo siguiente y compare sus resultados con los del estudio de caso. Analice en qué difieren y por qué.

- Calcule un intervalo de confianza del 99% de la media del diámetro.
- Calcule un intervalo de predicción del 99% en la medición del siguiente diámetro.
- Calcule un intervalo de tolerancia del 99% para la cobertura del 95% central de la distribución de diámetros.

9.28 En la sección 9.3 destacamos el concepto del "estimador más eficaz" comparando la varianza de dos estimadores insesgados $\hat{\theta}_1$ y $\hat{\theta}_2$. Sin embargo, esto no toma en cuenta el sesgo en el caso en que uno o ambos estimadores no son sesgados. Considere la cantidad

$$EME = E(\hat{\theta} - \theta),$$

donde EME denota el **error cuadrático medio**. El error cuadrático medio a menudo se utiliza para comparar dos estimadores $\hat{\theta}_1$ y $\hat{\theta}_2$ de θ , cuando uno o ambos no son sesgados porque i) es intuitivamente razonable y ii) se toma en cuenta para el sesgo. Demuestre que el EME se puede escribir como

$$\begin{aligned} EME &= E[\hat{\theta} - E(\hat{\theta})]^2 + [E(\hat{\theta}) - \theta]^2 \\ &= \text{Var}(\hat{\theta}) + [\text{sesgo}(\hat{\theta})]^2 \end{aligned}$$

9.29 Definamos $S^{*2} = \sum_{i=1}^n (X_i - \bar{X})^2 / n$. Demuestre que

$$E(S^{*2}) = [(n-1)/n]\sigma^2,$$

y, en consecuencia, que S^2 es un estimador sesgado para σ^2 .

9.29. Considere S^2 , el estimador de σ^2 , del ejercicio 9.29. Con frecuencia los analistas utilizan S^2 en lugar de dividir $\sum_{i=1}^n (X_i - \bar{X})^2$ entre $n - 1$, los grados de libertad en la muestra.

- a) ¿Cuál es el sesgo de S^2 ?
 b) Demuestre que el sesgo de S^2 se aproxima a cero a medida que $n \rightarrow \infty$.

9.31 Si X es una variable aleatoria binomial, demuestre que

- a) $\hat{p} = X/n$ es un estimador insesgado de p ;
 b) $p' = \frac{x + \sqrt{n}/2}{n + \sqrt{n}}$ es un estimador sesgado de p .

9.32 Demuestre que el estimador P' del ejercicio 9.31b) se vuelve no sesgado a medida que $n \rightarrow \infty$.

9.33 Compare S^2 y S'^2 (véase el ejercicio 9.29), los dos estimadores de σ^2 , para determinar cuál es más eficaz. Suponga que estos estimadores se obtienen usando X_1, X_2, \dots, X_n , las variables aleatorias independientes de $n(x; \mu, \sigma)$. ¿Cuál es el estimador más eficaz si se considera sólo la varianza de los estimadores? [Sugerencia: Utilice el teorema 8.4 y el hecho de que la varianza de χ^2_ν es 2ν , de la sección 6.7.]

9.34 Considere el ejercicio 9.33. Utilice el EME que se estudió en el ejercicio 9.28 para determinar qué estimador es más eficaz. Escriba

$$\frac{EME(S^2)}{EME(S'^2)}$$

9.8 Dos muestras: estimación de la diferencia entre dos medias

Si tenemos dos poblaciones con medias μ_1 y μ_2 , y varianzas σ_1^2 y σ_2^2 , respectivamente, el estadístico que da un estimador puntual de la diferencia entre μ_1 y μ_2 es $\bar{X}_1 - \bar{X}_2$. Por lo tanto, para obtener una estimación puntual de $\mu_1 - \mu_2$, se seleccionan dos muestras aleatorias independientes, una de cada población, de tamaños n_1 y n_2 , y se calcula $\bar{x}_1 - \bar{x}_2$, la diferencia de las medias muestrales. Evidentemente, debemos considerar la distribución muestral de $\bar{X}_1 - \bar{X}_2$.

De acuerdo con el teorema 8.3, podemos esperar que la distribución muestral de $\bar{X}_1 - \bar{X}_2$ esté distribuida de forma aproximadamente normal con media $\mu_{\bar{X}_1 - \bar{X}_2} = \mu_1 - \mu_2$ y desviación estándar $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}$. Por lo tanto, podemos asegurar, con una probabilidad de $1 - \alpha$, que la variable normal estándar

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$$

caerá entre $-z_{\alpha/2}$ y $z_{\alpha/2}$. Si nos remitimos una vez más a la figura 9.2, escribimos

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha.$$

Al sustituir para Z , establecemos de manera equivalente que

$$P\left(-z_{\alpha/2} < \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}} < z_{\alpha/2}\right) = 1 - \alpha,$$

que conduce al siguiente intervalo de confianza del $100(1 - \alpha)\%$ para $\mu_1 - \mu_2$.

Intervalo de confianza para $\mu_1 - \mu_2$ cuando se conocen σ_1^2 y σ_2^2

Si \bar{x}_1 y \bar{x}_2 son las medias de muestras aleatorias independientes de tamaños n_1 y n_2 , de poblaciones que tienen varianzas conocidas σ_1^2 y σ_2^2 , respectivamente, un intervalo de confianza del $100(1 - \alpha)\%$ para $\mu_1 - \mu_2$ es dado por

$$(\bar{x}_1 - \bar{x}_2) - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}},$$

donde $z_{\alpha/2}$ es el valor z que deja una área de $\alpha/2$ a la derecha.

El grado de confianza es exacto cuando las muestras se seleccionan de poblaciones normales. Para poblaciones no normales el teorema del límite central permite una buena aproximación para muestras de tamaño razonable.

Las condiciones experimentales y la unidad experimental

Para el caso en que se necesita estimar un intervalo de confianza sobre la diferencia entre dos medias se requiere considerar las condiciones experimentales durante el proceso de recolección de datos. Se supone que tenemos dos muestras aleatorias independientes de distribuciones con medias μ_1 y μ_2 , respectivamente. Es importante que las condiciones experimentales se parezcan al ideal descrito por las suposiciones tanto como sea posible. Con mucha frecuencia el experimentador debería planear la estrategia del experimento de acuerdo con esto. Para casi cualquier estudio de este tipo existe una *unidad experimental*, que es la parte del experimento que produce el error experimental y genera la varianza de la población que denominamos σ^2 . En un estudio farmacológico la unidad experimental es el paciente o el sujeto. En un experimento de agricultura puede ser una superficie de tierra. En un experimento químico puede ser una cantidad de materias primas. Es importante que las diferencias entre tales unidades tengan un impacto mínimo sobre los resultados. El experimentador tendrá un grado de seguridad de que las unidades experimentales no sesgarán los resultados si las condiciones que definen a las dos poblaciones se *asignan al azar* a las unidades experimentales. En los siguientes capítulos acerca de la prueba de hipótesis nos volveremos a concentrar en la aleatorización.

Ejemplo 9.10: Se llevó a cabo un experimento donde se compararon dos tipos de motores, el A y el B. Se midió el rendimiento de combustible en millas por galón. Se realizaron 50 experimentos con el motor tipo A y 75 con el motor tipo B. La gasolina utilizada y las demás condiciones se mantuvieron constantes. El rendimiento promedio de gasolina para el motor A fue de 36 millas por galón y el promedio para el motor B fue de 42 millas por galón. Calcule un intervalo de confianza del 96% sobre $\mu_B - \mu_A$, donde μ_A y μ_B corresponden a la media de la población del rendimiento de millas por galón para los motores A y B, respectivamente. Suponga que las desviaciones estándar de la población son 6 y 8 para los motores A y B, respectivamente.

Solución: La estimación puntual de $\mu_B - \mu_A$ es $\bar{x}_B - \bar{x}_A = 42 - 36 = 6$. Si usamos $\alpha = 0.04$, obtenemos $z_{0.02} = 2.05$ de la tabla A.3. Por lo tanto, sustituyendo en la fórmula anterior, el intervalo de confianza del 96% es

$$6 - 2.05\sqrt{\frac{64}{75} + \frac{36}{50}} < \mu_B - \mu_A < 6 + 2.05\sqrt{\frac{64}{75} + \frac{36}{50}},$$

o simplemente $3.43 < \mu_B - \mu_A < 8.57$. ▮

Este procedimiento para estimar la diferencia entre dos medias se aplica si se conocen σ_1^2 y σ_2^2 . Si las varianzas no se conocen y las dos distribuciones implicadas son aproximadamente normales, la distribución *t* resulta implicada como en el caso de una sola muestra. Si no se está dispuesto a suponer normalidad, muestras grandes (digamos mayores que 30) permitirán usar s_1 y s_2 en lugar de σ_1 y σ_2 , respectivamente, con el fundamento de que $s_1 \approx \sigma_1$ y $s_2 \approx \sigma_2$. De nuevo, por supuesto, el intervalo de confianza es aproximado.

Varianzas desconocidas pero iguales

Considere el caso donde se desconocen σ_1^2 y σ_2^2 . Si $\sigma_1^2 = \sigma_2^2 = \sigma^2$ obtenemos una variable normal estándar de la forma

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\sigma^2[(1/n_1) + (1/n_2)]}}$$

De acuerdo con el teorema 8.4, las dos variables aleatorias

$$\frac{(n_1 - 1)S_1^2}{\sigma^2} \quad \text{y} \quad \frac{(n_2 - 1)S_2^2}{\sigma^2}$$

tienen distribuciones chi cuadrada con $n_1 - 1$ y $n_2 - 1$ grados de libertad, respectivamente. Además, son variables chi cuadrada independientes, ya que las muestras aleatorias se seleccionaron de forma independiente. En consecuencia, su suma

$$V = \frac{(n_1 - 1)S_1^2}{\sigma^2} + \frac{(n_2 - 1)S_2^2}{\sigma^2} = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{\sigma^2}$$

tiene una distribución chi cuadrada con $v = n_1 + n_2 - 2$ grados de libertad.

Como se puede demostrar que las expresiones anteriores para Z y V son independientes, del teorema 8.5 se sigue que el estadístico

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\sigma^2[(1/n_1) + (1/n_2)]}} \bigg/ \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{\sigma^2(n_1 + n_2 - 2)}}$$

tiene la distribución t con $v = n_1 + n_2 - 2$ grados de libertad.

Se puede obtener una estimación puntual de la varianza común desconocida σ^2 agrupando las varianzas muestrales. Si representamos con S_p^2 al estimador agrupado, obtenemos lo siguiente,

Estimado
agrupado
de la varianza

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

Al sustituir S_p^2 en el estadístico T , obtenemos la forma menos engorrosa:

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{(1/n_1) + (1/n_2)}}$$

Si usamos el estadístico T , tenemos

$$P(-t_{\alpha/2} < T < t_{\alpha/2}) = 1 - \alpha,$$

donde $t_{\alpha/2}$ es el valor t con $n_1 + n_2 - 2$ grados de libertad, por arriba del cual encontramos una área de $\alpha/2$. Al sustituir por T en la desigualdad, escribimos

$$P \left[-t_{\alpha/2} < \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{(1/n_1) + (1/n_2)}} < t_{\alpha/2} \right] = 1 - \alpha.$$

Después de realizar las manipulaciones matemáticas de costumbre, se calculan la diferencia de las medias muestrales $\bar{x}_1 - \bar{x}_2$ y la varianza agrupada, y se obtiene el siguiente intervalo de confianza del $100(1 - \alpha)\%$ para $\mu_1 - \mu_2$.

Se observa con facilidad que el valor de s_p^2 es un promedio ponderado de las dos varianzas muestrales s_1^2 y s_2^2 , donde los pesos son los grados de libertad.

Intervalo de confianza para $\mu_1 - \mu_2$, $\sigma_1^2 = \sigma_2^2$ cuando se desconocen ambas varianzas

Si \bar{x}_1 y \bar{x}_2 son las medias de muestras aleatorias independientes con tamaños n_1 y n_2 , respectivamente, tomadas de poblaciones más o menos normales con varianzas iguales pero desconocidas, un intervalo de confianza del $100(1 - \alpha)\%$ para $\mu_1 - \mu_2$ es dado por

$$(\bar{x}_1 - \bar{x}_2) - t_{\alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + t_{\alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}},$$

donde s_p es la estimación agrupada de la desviación estándar de la población y $t_{\alpha/2}$ es el valor t con $v = n_1 + n_2 - 2$ grados de libertad, que deja una área de $\alpha/2$ a la derecha.

Ejemplo 9.11: En el artículo "Estructura comunitaria de los macroinvertebrados como un indicador de la contaminación de minas ácidas", publicado en el *Journal of Environmental Pollution*, se informa sobre una investigación realizada en Cane Creek, Alabama, para determinar la relación entre parámetros fisicoquímicos seleccionados y diversas mediciones de la estructura de la comunidad de macroinvertebrados. Una faceta de la investigación consistió en evaluar la efectividad de un índice numérico de la diversidad de especies para indicar la degradación del agua debida al desagüe ácido de una mina. Conceptualmente, un índice elevado de la diversidad de especies macroinvertebradas debería indicar un sistema acuático no contaminado; mientras que un índice bajo de esta diversidad indicaría un sistema acuático contaminado.

Se eligieron 2 estaciones de muestreo independientes para este estudio: una que se localiza corriente abajo del punto de descarga ácida de la mina y la otra ubicada corriente arriba. Para 12 muestras mensuales reunidas en la estación corriente abajo el índice de diversidad de especies tuvo un valor medio de $\bar{x}_1 = 3.11$ y una desviación estándar de $s_1 = 0.771$; mientras que 10 muestras reunidas mensualmente en la estación corriente arriba tuvieron un valor medio del índice $\bar{x}_2 = 2.04$ y una desviación estándar de $s_2 = 0.448$. Calculemos un intervalo de confianza del 90% para la diferencia entre las medias de la población de los dos sitios, suponiendo que las poblaciones se distribuyen de forma aproximadamente normal y que tienen varianzas iguales.

Solución: Representemos con μ_1 y μ_2 las medias de la población para los índices de diversidad de especies en las estaciones corriente abajo y corriente arriba, respectivamente. Deseamos encontrar un intervalo de confianza del 90% para $\mu_1 - \mu_2$. La estimación puntual de $\mu_1 - \mu_2$ es

$$\bar{x}_1 - \bar{x}_2 = 3.11 - 2.04 = 1.07$$

El estimado agrupado, s_p^2 , de la varianza común, σ^2 , es

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} = \frac{(11)(0.771^2) + (9)(0.448^2)}{12 + 10 - 2} = 0.417.$$

Al sacar la raíz cuadrada obtenemos $s_p = 0.646$. Si usamos $\alpha = 0.1$, encontramos en la tabla A.4 que $t_{0.05} = 1.725$ para $v = n_1 + n_2 - 2 = 20$ grados de libertad. Por lo tanto, el intervalo de confianza del 90% para $\mu_1 - \mu_2$ es

$$1.07 - (1.725)(0.646) \sqrt{\frac{1}{12} + \frac{1}{10}} < \mu_1 - \mu_2 < 1.07 + (1.725)(0.646) \sqrt{\frac{1}{12} + \frac{1}{10}},$$

que se simplifica a $0.593 < \mu_1 - \mu_2 < 1.547$.

Interpretación del intervalo de confianza

Para el caso de un solo parámetro el intervalo de confianza simplemente produce límites de error del parámetro. Los valores contenidos en el intervalo se deberían ver como valores razonables, dados los datos experimentales. En el caso de una diferencia entre dos medias, la interpretación se puede extender a una comparación de las dos medias. Por ejemplo, si tenemos gran confianza en que una diferencia $\mu_1 - \mu_2$ es positiva, sin duda inferiremos que $\mu_1 > \mu_2$ con poco riesgo de incurrir en un error. Así, en el ejemplo 9.11 tenemos un 90% de confianza en que el intervalo de 0.593 a 1.547 contiene la diferencia de las medias de la población para valores del índice de diversidad de especies en las dos estaciones. El hecho de que ambos límites de confianza sean positivos indica que, en promedio, el índice para la estación que se localiza corriente abajo del punto de descarga es mayor que el índice para la estación que se localiza corriente arriba.

Muestras de tamaños iguales

El procedimiento para construir intervalos de confianza para $\mu_1 - \mu_2$ cuando $\sigma_1 = \sigma_2 = \sigma$ pero ésta se desconoce, requiere suponer que las poblaciones son normales. Desviaciones ligeras de la suposición de varianzas iguales o de normalidad no alteran seriamente el grado de confianza en nuestro intervalo. (En el capítulo 10 se estudia un procedimiento para probar la igualdad de dos varianzas poblacionales desconocidas con base en la información que proporcionan las varianzas muestrales). Si las varianzas de la población son considerablemente diferentes, aún obtenemos resultados razonables cuando las poblaciones son normales, siempre y cuando $n_1 = n_2$. Por lo tanto, al planear un experimento se debería hacer un esfuerzo por igualar el tamaño de las muestras.

Varianzas desconocidas y distintas

Consideremos ahora el problema de calcular el estimado de un intervalo de $\mu_1 - \mu_2$ cuando no es probable que las varianzas de la población desconocidas sean iguales. El estadístico que se utiliza con mayor frecuencia en este caso es

$$T' = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{(S_1^2/n_1) + (S_2^2/n_2)}},$$

que tiene aproximadamente una distribución t con v grados de libertad, donde

$$v = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{[(s_1^2/n_1)^2/(n_1 - 1)] + [(s_2^2/n_2)^2/(n_2 - 1)]}.$$

Como v rara vez es un entero, lo *redondeamos* al número entero menor más cercano. El estimado anterior de los grados de libertad se denomina aproximación de Satterthwaite (Satterthwaite, 1946, en la bibliografía).

Con el estadístico T' , escribimos

$$P(-t_{\alpha/2} < T' < t_{\alpha/2}) \approx 1 - \alpha,$$

donde $t_{\alpha/2}$ es el valor de la distribución t con v grados de libertad por arriba del cual encontramos una área de $\alpha/2$. Al sustituir para T' en la desigualdad y seguir los mismos pasos que antes, establecemos el resultado final.

Intervalo de confianza para $\mu_1 - \mu_2$, $\sigma_1^2 \neq \sigma_2^2$ y ambas varianzas se desconocen

Si \bar{x}_1 y s_1^2 y \bar{x}_2 y s_2^2 son las medias y varianzas de muestras aleatorias independientes de tamaños n_1 y n_2 , respectivamente, tomadas de poblaciones aproximadamente normales con varianzas desconocidas y diferentes, un intervalo de confianza aproximado del 100(1 - α)% para $\mu_1 - \mu_2$ es dado por

$$(\bar{x}_1 - \bar{x}_2) - t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}},$$

donde $t_{\alpha/2}$ es el valor t con

$$v = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{[(s_1^2/n_1)^2/(n_1 - 1)] + [(s_2^2/n_2)^2/(n_2 - 1)]}$$

grados de libertad, que deja una área de $\alpha/2$ a la derecha.

Observe que la expresión para el valor v anterior incluye variables aleatorias y, por consiguiente, v es un *estimado* de los grados de libertad. En las aplicaciones este estimado no será un número entero, de manera que el analista lo debe redondear al entero menor más cercano para lograr la confianza que se busca.

Antes de ilustrar el intervalo de confianza anterior con un ejemplo deberíamos señalar que todos los intervalos de confianza para $\mu_1 - \mu_2$ tienen la misma forma general, como los de una sola media; a saber, se pueden escribir como

$$\text{estimación puntual} \pm t_{\alpha/2} \widehat{\text{e.e.}} (\text{estimación puntual})$$

o

$$\text{estimación puntual} \pm z_{\alpha/2} \text{e.e.} (\text{estimación puntual}).$$

Por ejemplo, en el caso donde $\sigma_1 = \sigma_2 = \sigma$, el error estándar estimado de $\bar{x}_1 - \bar{x}_2$ es $s_p \sqrt{1/n_1 + 1/n_2}$. Para el caso donde $\sigma_1^2 \neq \sigma_2^2$,

$$\widehat{\text{e.e.}}(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$$

Ejemplo 9.12: El Departamento de zoología de Virginia Tech llevó a cabo un estudio para estimar la diferencia en la cantidad de ortofósforo químico medido en dos estaciones diferentes del río James. El ortofósforo se mide en miligramos por litro. Se reunieron 15 muestras de la estación 1 y 12 muestras de la estación 2. Las 15 muestras de la estación 1 tuvieron un contenido promedio de ortofósforo de 3.84 miligramos por litro y una desviación estándar de 3.07 miligramos por litro; en tanto que las 12 muestras de la estación 2 tuvieron un contenido promedio de 1.49 miligramos por litro y una desviación estándar de 0.80 miligramos por litro. Calcule un intervalo de confianza de 95% para la diferencia en el contenido promedio verdadero de ortofósforo en estas dos estaciones. Suponga que las observaciones provienen de poblaciones normales con varianzas diferentes.

Solución: Para la estación 1 tenemos $\bar{x}_1 = 3.84$, $s_1 = 3.07$ y $n_1 = 15$. Para la estación 2, $\bar{x}_2 = 1.49$, $s_2 = 0.80$ y $n_2 = 12$. Queremos obtener un intervalo de confianza del 95% para $\mu_1 - \mu_2$.

Como se suponen varianzas de la población diferentes, sólo podemos calcular un intervalo de confianza aproximado del 95% basado en la distribución t con ν grados de libertad, donde

$$\nu = \frac{(3.07^2/15 + 0.80^2/12)^2}{[(3.07^2/15)^2/14] + [(0.80^2/12)^2/11]} = 16.3 \approx 16.$$

Nuestra estimación puntual de $\mu_1 - \mu_2$ es

$$\bar{x}_1 - \bar{x}_2 = 3.84 - 1.49 = 2.35.$$

Si usamos $\alpha = 0.05$, en la tabla A.4 encontramos que $t_{0.025} = 2.120$ para $\nu = 16$ grados de libertad. Por lo tanto, el intervalo de confianza del 95% para $\mu_1 - \mu_2$ es

$$2.35 - 2.120\sqrt{\frac{3.07^2}{15} + \frac{0.80^2}{12}} < \mu_1 - \mu_2 < 2.35 + 2.120\sqrt{\frac{3.07^2}{15} + \frac{0.80^2}{12}},$$

que se simplifica a $0.60 < \mu_1 - \mu_2 < 4.10$. En consecuencia, tenemos un 95% de confianza en que el intervalo de 0.60 a 4.10 miligramos por litro contiene la diferencia del promedio verdadero del ortofósforo que contienen estos dos lugares. ■

Cuando se desconocen dos varianzas de la población, la suposición de varianzas iguales o diferentes podría ser precaria. En la sección 10.10 se presentará un procedimiento que ayudará a distinguir entre las situaciones con la misma varianza y con varianza diferente.

9.9 Observaciones pareadas

Ahora estudiaremos los procedimientos de estimación para la diferencia de dos medias cuando las muestras no son independientes y las varianzas de las dos poblaciones no son necesariamente iguales. La situación que se considera aquí tiene que ver con una condición experimental muy especial, a saber, *las observaciones pareadas*. A diferencia de la situación que se describió antes, las condiciones de las dos poblaciones no se asignan de forma aleatoria a las unidades experimentales. Más bien, cada unidad experimental homogénea recibe ambas condiciones de la población; como resultado, cada unidad experimental tiene un par de observaciones, una para cada población. Por ejemplo, si realizamos una prueba de una nueva dieta con 15 individuos, los pesos antes y después de seguir la dieta conforman la información de las dos muestras. Las dos poblaciones son “antes” y “después”, y la unidad experimental es el individuo. Evidentemente, las observaciones en un par tienen algo en común. Para determinar si la dieta es efectiva consideramos las diferencias d_1, d_2, \dots, d_n en las observaciones pareadas. Estas diferencias son los valores de una muestra aleatoria D_1, D_2, \dots, D_n de una población de diferencias, que supondremos distribuidas normalmente, con media $\mu_D = \mu_1 - \mu_2$ y varianza σ_D^2 . Estimamos σ_D^2 mediante s_D^2 , la varianza de las diferencias que constituyen nuestra muestra. El estimador puntual de μ_D es dado por \bar{D} .

¿Cuándo debe hacerse el pareado?

Parear observaciones en un experimento es una estrategia que se puede emplear en muchos campos de aplicación. Se expondrá al lector a tal concepto en el material relacionado con

la prueba de hipótesis en el capítulo 10 y en los temas de diseño experimental en los capítulos 13 y 15. Al seleccionar unidades experimentales relativamente homogéneas (dentro de las unidades) y permitir que cada unidad experimente ambas condiciones de la población, se reduce la varianza del error experimental efectiva (en este caso σ_D^2). El lector puede visualizar la i -ésima diferencia del par como

$$D_i = X_{1i} - X_{2i}.$$

Como las dos observaciones se toman de la unidad experimental de la muestra no son independientes y, de hecho,

$$\text{Var}(D_i) = \text{Var}(X_{1i} - X_{2i}) = \sigma_1^2 + \sigma_2^2 - 2 \text{Cov}(X_{1i}, X_{2i}).$$

Entonces, de manera intuitiva, se espera que σ_D^2 debería reducirse debido a la similitud en la naturaleza de los "errores" de las dos observaciones dentro de una unidad experimental, a lo cual se llega mediante la expresión anterior. En realidad se espera que, si la unidad es homogénea, la covarianza sea positiva. Como resultado, la ganancia en calidad del intervalo de confianza sobre la que se obtuvo sin parear es mayor cuando hay homogeneidad dentro de las unidades y cuando las diferencias grandes van de una a otra unidad. Se debería tener en cuenta que el desempeño del intervalo de confianza dependerá del error estándar de \bar{D} , que es, por supuesto, σ_D/\sqrt{n} , donde n es el número de pares. Como indicamos antes, la intención al parear es reducir σ_D .

Equilibrio entre reducir la varianza y perder grados de libertad

Al comparar los intervalos de confianza obtenidos con y sin pareado es evidente que hay un intercambio implicado. Aunque en realidad el pareado debería reducir la varianza y, por lo tanto, el error estándar de la estimación puntual, los grados de libertad disminuyen al reducir el problema a uno con una sola muestra. Como resultado, el punto $t_{\alpha/2}$ ligado al error estándar se ajusta en concordancia. De esta manera, el pareado podría resultar contraproducente. Esto ocurriría con certeza si se experimenta sólo una reducción modesta en la varianza (a través de σ_D^2) mediante el pareado.

Otra ilustración del pareado implicaría elegir n pares de sujetos, donde cada par tenga una característica similar, como el coeficiente intelectual (CI), la edad o la raza, y luego para cada par seleccionar un miembro al azar para obtener un valor de X_1 , dejando que el otro miembro proporcione el valor de X_2 . En este caso, X_1 y X_2 podrían representar las calificaciones obtenidas por dos individuos con igual CI cuando uno es asignado al azar a un grupo que usa el método de enseñanza convencional y al otro a un grupo que utiliza materiales programados.

Se puede establecer un intervalo de confianza del $100(1 - \alpha)\%$ para μ_D escribiendo

$$P(-t_{\alpha/2} < T < t_{\alpha/2}) = 1 - \alpha,$$

donde $T = \frac{\bar{D} - \mu_D}{S_D/\sqrt{n}}$ y $t_{\alpha/2}$, como antes, es un valor de la distribución t con $n - 1$ grados de libertad.

En la actualidad se acostumbra reemplazar T por su definición en la desigualdad anterior y desarrollar los pasos matemáticos que conduzcan al siguiente intervalo de confianza del $100(1 - \alpha)\%$ para $\mu_1 - \mu_2 = \mu_D$.

Intervalo de confianza para $\mu_D = \mu_1 - \mu_2$ para observaciones pareadas

Si \bar{d} y s_d son la media y la desviación estándar, respectivamente, de las diferencias distribuidas normalmente de n pares aleatorios de mediciones, un intervalo de confianza del $100(1 - \alpha)\%$ para $\mu_D = \mu_1 - \mu_2$ es

$$\bar{d} - t_{\alpha/2} \frac{s_d}{\sqrt{n}} < \mu_D < \bar{d} + t_{\alpha/2} \frac{s_d}{\sqrt{n}},$$

donde $t_{\alpha/2}$ es el valor t con $v = n - 1$ grados de libertad, que deja una área de $\alpha/2$ a la derecha.

Ejemplo 9.13: Un estudio publicado en *Chemosphere* reporta los niveles de la dioxina TCDD en 20 veteranos de Vietnam de Massachusetts, quienes posiblemente estuvieron expuestos al agente naranja. En la tabla 9.1 se presentan los niveles de TCDD en plasma y tejido adiposo.

Calcule un intervalo de confianza del 95% para $\mu_1 - \mu_2$, donde μ_1 y μ_2 representen las medias verdaderas de los niveles de TCDD en plasma y en tejido adiposo, respectivamente. Suponga que la distribución de las diferencias es casi normal.

Tabla 9.1: Datos para el ejemplo 9.13.

Veterano	Niveles de TCDD en plasma	Niveles de TCDD en tejido adiposo	d_i	Veterano	Niveles de TCDD en plasma	Niveles de TCDD en tejido adiposo	d_i
1	2.5	4.9	-2.4	11	6.9	7.0	-0.1
2	3.1	5.9	-2.8	12	3.3	2.9	0.4
3	2.1	4.4	-2.3	13	4.6	4.6	0.0
4	3.5	6.9	-3.4	14	1.6	1.4	0.2
5	3.1	7.0	-3.9	15	7.2	7.7	-0.5
6	1.8	4.2	-2.4	16	1.8	1.1	0.7
7	6.0	10.0	-4.0	17	20.0	11.0	9.0
8	3.0	5.5	-2.5	18	2.0	2.5	-0.5
9	36.0	41.0	-5.0	19	2.5	2.3	0.2
10	4.7	4.4	0.3	20	4.1	2.5	1.6

Reproducido de *Chemosphere*, Vol. 20, Núms. 7-9 (tablas I y II), Schecter *et al.*, "Partitioning 2, 3, 7, 8-chlorinated dibenzo-p-dioxins and dibenzofurans between adipose tissue and plasma lipid of 20 Massachusetts Vietnam veterans", pp. 954-955, Derechos reservados ©1990, con autorización de Elsevier.

Solución: Buscamos un intervalo de confianza del 95% para $\mu_1 - \mu_2$. Como las observaciones están pareadas, $\mu_1 - \mu_2 = \mu_D$. La estimación puntual de μ_D es $\bar{d} = -0.87$. La desviación estándar s_d de las diferencias muestrales es

$$s_d = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (d_i - \bar{d})^2} = \sqrt{\frac{168.4220}{19}} = 2.9773.$$

Si usamos $\alpha = 0.05$, en la tabla A.4 encontramos que $t_{0.025} = 2.093$ para $v = n - 1 = 19$ grados de libertad. Por lo tanto, el intervalo de confianza del 95% es

$$-0.8700 - (2.093) \left(\frac{2.9773}{\sqrt{20}} \right) < \mu_D < -0.8700 + (2.093) \left(\frac{2.9773}{\sqrt{20}} \right).$$

o simplemente $-2.2634 < \mu_D < 0.5234$, de lo cual concluimos que no hay diferencia significativa entre el nivel medio de TCDD en plasma y el nivel medio de TCDD en tejido adiposo.

Ejercicios

9.35 Una muestra aleatoria de tamaño $n_1 = 25$, tomada de una población normal con una desviación estándar $\sigma_1 = 5$, tiene una media $\bar{x}_1 = 80$. Una segunda muestra aleatoria de tamaño $n_2 = 36$, que se toma de una población normal diferente con una desviación estándar $\sigma_2 = 3$, tiene una media $\bar{x}_2 = 75$. Calcule un intervalo de confianza del 94% para $\mu_1 - \mu_2$.

9.36 Se comparan las resistencias de dos clases de hilo. Se prueban 50 piezas de cada clase de hilo en condiciones similares. La marca A tiene una resistencia a la tensión promedio de 78.3 kilogramos, con una desviación estándar de 5.6 kilogramos; en tanto que la marca B tiene una resistencia a la tensión promedio de 87.2 kilogramos con una desviación estándar de 6.3 kilogramos. Construya un intervalo de confianza del 95% para la diferencia de las medias de la población.

9.37 Se realiza un estudio para determinar si cierto tratamiento tiene algún efecto sobre la cantidad de metal que se elimina en una operación de encurtido. Una muestra aleatoria de 100 piezas se sumerge en un baño por 24 horas sin el tratamiento, lo que produce un promedio de 12.2 milímetros de metal eliminados y una desviación estándar muestral de 1.1 milímetros. Una segunda muestra de 200 piezas se somete al tratamiento, seguido de 24 horas de inmersión en el baño, lo que da como resultado una eliminación promedio de 9.1 milímetros de metal, con una desviación estándar muestral de 0.9 milímetros. Calcule un estimado del intervalo de confianza del 98% para la diferencia entre las medias de las poblaciones. ¿El tratamiento parece reducir la cantidad media del metal eliminado?

9.38 En un proceso químico por lotes se comparan los efectos de dos catalizadores sobre la potencia de la reacción del proceso. Se prepara una muestra de 12 lotes utilizando el catalizador 1 y una muestra de 10 lotes utilizando el catalizador 2. Los 12 lotes para los que se utilizó el catalizador 1 en la reacción dieron un rendimiento promedio de 85 con una desviación estándar muestral de 4; en tanto que para la segunda muestra, la de 10 lotes, el promedio fue de 81, con una desviación estándar muestral de 5. Calcule un intervalo de confianza del 90% para la diferencia entre las medias de la población, suponiendo que las poblaciones se distribuyen de forma aproximadamente normal y que tienen varianzas iguales.

9.39 Los estudiantes pueden elegir entre un curso de física de tres semestres-hora sin laboratorio y un curso de cuatro semestres-hora con laboratorio. El examen

final escrito es el mismo para ambos cursos. Si 12 estudiantes del curso con laboratorio obtienen una calificación promedio de 84, con una desviación estándar de 4, y 18 estudiantes del grupo sin laboratorio obtienen una calificación promedio de 77, con una desviación estándar de 6, calcule un intervalo de confianza del 99% para la diferencia entre las calificaciones promedio para ambos cursos. Suponga que las poblaciones se distribuyen de forma aproximadamente normal y que tienen varianzas iguales.

9.40 En un estudio que se lleva a cabo en Virginia Tech sobre el desarrollo de micorriza, una relación simbiótica entre las raíces de árboles y un hongo, en la cual se transfieren minerales del hongo a los árboles y azúcares de los árboles a los hongos, se cultivaron en un invernadero 20 robles rojos que fueron expuestos al hongo *Pisolithus tinctorius*. Todos los árboles se plantaron en el mismo tipo de suelo y recibieron la misma cantidad de luz solar y agua. La mitad no recibió nitrógeno en el momento de plantarlos y sirvió como control, y la otra mitad recibió 368 ppm de nitrógeno en forma de NaNO_3 . Después de 140 días se registraron los siguientes pesos de los tallos, en gramos:

Sin nitrógeno	Con nitrógeno
0.32	0.26
0.53	0.43
0.28	0.47
0.37	0.49
0.47	0.52
0.43	0.75
0.36	0.79
0.42	0.86
0.38	0.62
0.43	0.46

Construya un intervalo de confianza del 95% para la diferencia entre los pesos medios de los tallos que no recibieron nitrógeno y los que recibieron 368 ppm de nitrógeno. Suponga que las poblaciones están distribuidas normalmente y que tienen varianzas iguales.

9.41 Los siguientes datos representan el tiempo, en días, que pacientes tratados al azar con uno de dos medicamentos para curar infecciones graves de la vejiga tardaron en recuperarse:

Medicamento 1	Medicamento 2
$n_1 = 14$	$n_2 = 16$
$\bar{x}_1 = 17$	$\bar{x}_2 = 19$
$s_1^2 = 1.5$	$s_2^2 = 1.8$

Calcule un intervalo de confianza del 99% para la diferencia $\mu_2 - \mu_1$ en los tiempos medios de recuperación para los dos medicamentos. Suponga poblaciones normales que tienen varianzas iguales.

9.42 Un experimento publicado en *Popular Science* comparó el ahorro de combustible para dos tipos de camiones compactos que funcionan con diesel y están equipados de forma similar. Suponga que se utilizaron 12 camiones Volkswagen y 10 Toyota en pruebas con una velocidad constante de 90 kilómetros por hora. Si los 12 camiones Volkswagen promedian 16 kilómetros por litro con una desviación estándar de 1.0 kilómetros por litro, y los 10 Toyota promedian 11 kilómetros por litro con una desviación estándar de 0.8 kilómetros por litro, construya un intervalo de confianza del 90% para la diferencia entre los kilómetros promedio por litro de estos dos camiones compactos. Suponga que las distancias por litro para cada modelo de camión están distribuidas de forma aproximadamente normal y que tienen varianzas iguales.

9.43 Una empresa de taxis trata de decidir si comprará neumáticos de la marca *A* o de la marca *B* para su flotilla de taxis. Para estimar la diferencia entre las dos marcas realiza un experimento utilizando 12 neumáticos de cada marca, los cuales utiliza hasta que se desgastan. Los resultados son:

Marca *A*: $\bar{x}_1 = 36,300$ kilómetros,

$s_1 = 5000$ kilómetros.

Marca *B*: $\bar{x}_2 = 38,100$ kilómetros,

$s_2 = 6100$ kilómetros.

Calcule un intervalo de confianza del 95% para $\mu_A - \mu_B$, suponiendo que las poblaciones se distribuyen de forma aproximadamente normal. Puede no suponer que las varianzas son iguales.

9.44 Con referencia al ejercicio 9.43, calcule un intervalo de confianza del 99% para $\mu_1 - \mu_2$ si se asignan al azar neumáticos de las dos marcas a las ruedas traseras izquierda y derecha de 8 taxis y se registran las siguientes distancias, en kilómetros:

Taxi	Marca A	Marca B
1	34,400	36,700
2	45,500	46,800
3	36,700	37,700
4	32,000	31,100
5	48,400	47,800
6	32,800	36,400
7	38,100	38,900
8	30,100	31,500

Suponga que las diferencias de las distancias se distribuyen de forma aproximadamente normal.

9.45 El gobierno otorgó fondos para los departamentos de agricultura de 9 universidades para probar las

capacidades de cosecha de dos nuevas variedades de trigo. Cada variedad se siembra en parcelas con la misma área en cada universidad, y las cosechas, en kilogramos por parcela, se registran como sigue:

Variedad	Universidad								
	1	2	3	4	5	6	7	8	9
1	38	23	35	41	44	29	37	31	38
2	45	25	31	38	50	33	36	40	43

Calcule un intervalo de confianza del 95% para la diferencia media entre las cosechas de las dos variedades, suponiendo que las diferencias entre las cosechas se distribuyen de forma aproximadamente normal. Explique por qué es necesario el pareado en este problema.

9.46 Los siguientes datos representan el tiempo de duración de películas producidas por dos empresas cinematográficas.

Empresa	Tiempo (minutos)								
I	103	94	110	87	98				
II	97	82	123	92	175	88	118		

Calcule un intervalo de confianza del 90% para la diferencia entre la duración promedio de las películas que producen las dos empresas. Suponga que las diferencias en la duración se distribuyen de forma aproximadamente normal y que tienen varianzas distintas.

9.47 La revista *Fortune* (marzo de 1997) publicó la rentabilidad total de los inversionistas durante los 10 años anteriores a 1996 y también la de 431 empresas en ese mismo año. A continuación se lista la rentabilidad total para 10 de las empresas. Calcule un intervalo de confianza del 95% para el cambio promedio en el porcentaje de rentabilidad de los inversionistas.

Empresa	Rentabilidad total para los inversionistas	
	1986-96	1996
Coca-Cola	29.8%	43.3%
Mirage Resorts	27.9%	25.4%
Merck	22.1%	24.0%
Microsoft	44.5%	88.3%
Johnson & Johnson	22.2%	18.1%
Intel	43.8%	131.2%
Pfizer	21.7%	34.0%
Procter & Gamble	21.9%	32.1%
Berkshire Hathaway	28.3%	6.2%
S&P 500	11.8%	20.3%

9.48 Una empresa automotriz está considerando dos tipos de baterías para sus vehículos. Con ese fin reúne información muestral sobre la vida de las baterías. Utiliza para ello 20 baterías del tipo *A* y 20 baterías del tipo *B*. El resumen de los estadísticos es $\bar{x}_A = 32.91$.

$\bar{x}_B = 30.47$, $s_A = 1.57$ y $s_B = 1.74$. Suponga que los datos de cada batería se distribuyen normalmente y que $\sigma_A = \sigma_B$.

- a) Calcule un intervalo de confianza del 95% para $\mu_A - \mu_B$.
- b) Del inciso a) saque algunas conclusiones que le ayuden a la empresa a decidir si debería utilizar la batería A o la B.

9.49 Se considera usar dos marcas diferentes de pintura vinílica. Se seleccionaron 15 especímenes de cada tipo de pintura, para los cuales los tiempos de secado en horas fueron los siguientes:

Pintura A					Pintura B				
3.5	2.7	3.9	4.2	3.6	4.7	3.9	4.5	5.5	4.0
2.7	3.3	5.2	4.2	2.9	5.3	4.3	6.0	5.2	3.7
4.4	5.2	4.0	4.1	3.4	5.5	6.2	5.1	5.4	4.8

Suponga que el tiempo de secado se distribuye normalmente, con $\sigma_A = \sigma_B$. Calcule un intervalo de confianza del 95% de $\mu_B - \mu_A$, donde μ_A y μ_B son los tiempos medios de secado.

9.50 A dos grupos de ratas diabéticas se les suministran dos niveles de dosis de insulina (alto y bajo) para verificar la capacidad de fijación de esta hormona. Se obtuvieron los siguientes datos.

Dosis baja: $n_1 = 8$ $\bar{x}_1 = 1.98$ $s_1 = 0.51$
 Dosis alta: $n_2 = 13$ $\bar{x}_2 = 1.30$ $s_2 = 0.35$

Suponga que las varianzas son iguales. Determine un intervalo de confianza del 95% para la diferencia en la capacidad promedio verdadera de fijación de la insulina entre las dos muestras.

9.10 Una sola muestra: estimación de una proporción

El estadístico $\hat{P} = X/n$, en donde X representa el número de éxitos en n ensayos, provee un estimador puntual de la proporción p en un experimento binomial. Por lo tanto, la proporción de la muestra $\hat{p} = x/n$ se utilizará como el estimador puntual del parámetro p .

Si no se espera que la proporción p desconocida esté demasiado cerca de 0 o de 1, se puede establecer un intervalo de confianza para p considerando la distribución muestral de \hat{P} . Si en cada ensayo binomial asignamos el valor 0 a un fracaso y el valor 1 a un éxito, el número de éxitos, x , se puede interpretar como la suma de n valores que consta sólo de ceros y unos, y \hat{p} es sólo la media muestral de esos n valores. En consecuencia, por el teorema del límite central, para n suficientemente grande \hat{P} está distribuida de forma casi normal con media

$$\mu_{\hat{p}} = E(\hat{P}) = E\left(\frac{X}{n}\right) = \frac{np}{n} = p$$

y varianza

$$\sigma_{\hat{p}}^2 = \sigma_{X/n}^2 = \frac{\sigma_X^2}{n^2} = \frac{npq}{n^2} = \frac{pq}{n}.$$

Por lo tanto, podemos afirmar que

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha, \text{ con } Z = \frac{\hat{P} - p}{\sqrt{pq/n}},$$

y $z_{\alpha/2}$ es el valor por arriba del cual encontramos una área de $\alpha/2$ debajo de la curva normal estándar. Al sustituir para Z escribimos

$$P\left(-z_{\alpha/2} < \frac{\hat{P} - p}{\sqrt{pq/n}} < z_{\alpha/2}\right) = 1 - \alpha.$$

Cuando n es grande se introduce un error muy pequeño sustituyendo el estimado puntual $\hat{p} = x/n$ para la p debajo del signo de radical. Entonces podemos escribir

$$P\left(\hat{P} - z_{\alpha/2}\sqrt{\frac{\hat{p}\hat{q}}{n}} < p < \hat{P} + z_{\alpha/2}\sqrt{\frac{\hat{p}\hat{q}}{n}}\right) \approx 1 - \alpha.$$

Por otro lado, al resolver para p en la desigualdad cuadrática anterior,

$$-z_{\alpha/2} < \frac{\hat{p} - p}{\sqrt{pq/n}} < z_{\alpha/2},$$

obtenemos otra forma del intervalo de confianza para p con los siguientes límites:

$$\frac{\hat{p} + \frac{z_{\alpha/2}^2}{2n}}{1 + \frac{z_{\alpha/2}^2}{n}} \pm \frac{z_{\alpha/2}}{1 + \frac{z_{\alpha/2}^2}{n}} \sqrt{\frac{\hat{p}\hat{q}}{n} + \frac{z_{\alpha/2}^2}{4n^2}}.$$

Para una muestra aleatoria de tamaño n se calcula la proporción muestral $\hat{p} = x/n$ y se pueden obtener los siguientes intervalos de confianza aproximados del $100(1 - \alpha)\%$ para p .

Intervalos de confianza para p de una muestra grande

Si \hat{p} es la proporción de éxitos en una muestra aleatoria de tamaño n , y $\hat{q} = 1 - \hat{p}$, un intervalo de confianza aproximado del $100(1 - \alpha)\%$ para el parámetro binomial p se obtiene por medio de (método 1)

$$\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}} < p < \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$$

o mediante (método 2)

$$\frac{\hat{p} + \frac{z_{\alpha/2}^2}{2n}}{1 + \frac{z_{\alpha/2}^2}{n}} - \frac{z_{\alpha/2}}{1 + \frac{z_{\alpha/2}^2}{n}} \sqrt{\frac{\hat{p}\hat{q}}{n} + \frac{z_{\alpha/2}^2}{4n^2}} < p < \frac{\hat{p} + \frac{z_{\alpha/2}^2}{2n}}{1 + \frac{z_{\alpha/2}^2}{n}} + \frac{z_{\alpha/2}}{1 + \frac{z_{\alpha/2}^2}{n}} \sqrt{\frac{\hat{p}\hat{q}}{n} + \frac{z_{\alpha/2}^2}{4n^2}}.$$

donde $z_{\alpha/2}$ es el valor z que deja una área de $\alpha/2$ a la derecha.

Cuando n es pequeña y se cree que la proporción desconocida p se acerca a 0 o a 1, el procedimiento del intervalo de confianza que se establece aquí no es confiable y, por lo tanto, no se debería emplear. Para estar seguros se requiere que tanto $n\hat{p}$ como $n\hat{q}$ sean mayores que o iguales a 5. Los métodos para calcular un intervalo de confianza para el parámetro binomial p también se pueden aplicar cuando se está utilizando la distribución binomial con el fin de aproximar la distribución hipergeométrica; es decir, cuando n es pequeña respecto a N , como se ilustra en el ejemplo 9.14.

Observe que, aunque el método 2 produce resultados más precisos, su cálculo es más complicado, y la ventaja en precisión que brinda disminuye cuando el tamaño de la muestra es lo suficientemente grande. Debido a esto en la práctica es más común utilizar el método 1.

Ejemplo 9.14: En una muestra aleatoria de $n = 500$ familias que tienen televisores en la ciudad de Hamilton, Canadá, se encuentra que $x = 340$ están suscritas a HBO. Calcule un intervalo de confianza del 95% para la proporción real de familias que tienen televisores en esta ciudad y están suscritas a HBO.

Solución: La estimación puntual de p es $\hat{p} = 340/500 = 0.68$. Si usamos la tabla A.3, encontramos que $z_{0.025} = 1.96$. Por lo tanto, si utilizamos el método 1, el intervalo de confianza del 95% para p es

$$0.68 - 1.96 \sqrt{\frac{(0.68)(0.32)}{500}} < p < 0.68 + 1.96 \sqrt{\frac{(0.68)(0.32)}{500}},$$

que se simplifica a $0.6391 < p < 0.7209$.

Si utilizamos el segundo método, obtenemos

$$\frac{0.68 + \frac{1.96^2}{(2)(500)}}{1 + \frac{1.96^2}{500}} \pm \frac{1.96}{1 + \frac{1.96^2}{500}} \sqrt{\frac{(0.68)(0.32)}{500} + \frac{1.96^2}{(4)(500^2)}} = 0.6786 \pm 0.0408,$$

que se simplifica a $0.6378 < p < 0.7194$. Aparentemente, cuando n es grande (500 en este caso) ambos métodos producen resultados muy similares. ▮

Si p es el valor central de un intervalo de confianza del $100(1 - \alpha)\%$, entonces \hat{p} estima p sin error. Sin embargo, la mayoría de las veces \hat{p} no será exactamente igual a p y el estimado puntual será erróneo. El tamaño de este error será la diferencia positiva que separa a p de \hat{p} , y podemos tener una confianza del $100(1 - \alpha)\%$ de que tal diferencia no excederá a $z_{\alpha/2} \sqrt{\hat{p}\hat{q}/n}$. Si dibujamos un diagrama de un intervalo de confianza típico, como el de la figura 9.6, podemos ver esto fácilmente. En este caso utilizamos el método 1 para estimar el error.

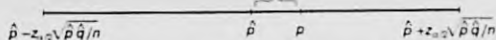


Figura 9.6: Error en la estimación de p por medio de \hat{p} .

Teorema 9.3: Si \hat{p} se utiliza como un estimado de p , podemos tener un $100(1 - \alpha)\%$ de confianza en que el error no excederá a $z_{\alpha/2} \sqrt{\hat{p}\hat{q}/n}$.

En el ejemplo 9.14 tenemos un 95% de confianza en que la proporción de la muestra $\hat{p} = 0.68$ difiere de la verdadera proporción p en una cantidad que no excede a 0.04.

Selección del tamaño de la muestra

Determinemos ahora qué tan grande debe ser una muestra para poder estar seguros de que el error al estimar p será menor que una cantidad específica e . Por medio del teorema 9.3, debemos elegir una n tal que $z_{\alpha/2} \sqrt{\hat{p}\hat{q}/n} = e$.

Teorema 9.4: Si \hat{p} se utiliza como un estimado de p , podemos tener un $100(1 - \alpha)\%$ de confianza en que el error será menor que una cantidad específica e cuando el tamaño de la muestra sea aproximadamente

$$n = \frac{z_{\alpha/2}^2 \hat{p}\hat{q}}{e^2}.$$

El teorema 9.4 es algo engañoso, pues debemos utilizar \hat{p} para determinar el tamaño n de la muestra, pero \hat{p} se calcula a partir de la muestra. Si se puede hacer una estimación burda de p sin tomar una muestra, se podría usar este valor para determinar n . A falta de tal estimado, podríamos tomar una muestra preliminar de tamaño $n \geq 30$ para proporcionar un estimado de p . Si utilizamos el teorema 9.4 podríamos determinar aproximadamente cuántas observaciones se necesitan para proporcionar el grado de precisión deseado. Observe que los valores fraccionarios de n se redondean al siguiente número entero mayor.

Ejemplo 9.15: ¿Qué tan grande debe ser una muestra en el ejemplo 9.14 si queremos tener un 95% de confianza en que la estimación de p esté dentro de 0.02 del valor verdadero?

Solución: Tratemos a las 500 familias como una muestra preliminar que proporciona una estimación $\hat{p} = 0.68$. Entonces, mediante el teorema 9.4,

$$n = \frac{(1.96)^2(0.68)(0.32)}{(0.02)^2} = 2089.8 \approx 2090.$$

Por lo tanto, si basamos nuestra estimación de p en una muestra aleatoria de tamaño 2090, podemos tener un 95% de confianza en que nuestra proporción muestral no diferirá de la proporción verdadera en más de 0.02. ■

Ocasionalmente será poco práctico obtener una estimación de p que se utilice para determinar el tamaño muestral para un grado específico de confianza. Si esto sucede, se establece un límite superior para n al notar que $\hat{p}\hat{q} = \hat{p}(1 - \hat{p})$, que debe ser a lo sumo $1/4$, ya que \hat{p} debe caer entre 0 y 1. Este hecho se verifica completando el cuadrado. Por consiguiente,

$$\hat{p}(1 - \hat{p}) = -(\hat{p}^2 - \hat{p}) = \frac{1}{4} - \left(\hat{p}^2 - \hat{p} + \frac{1}{4}\right) = \frac{1}{4} - \left(\hat{p} - \frac{1}{2}\right)^2,$$

que siempre es menor que $1/4$ excepto cuando $\hat{p} = 1/2$ y entonces $\hat{p}\hat{q} = 1/4$. Por lo tanto, si sustituimos $\hat{p} = 1/2$ en la fórmula para n del teorema 9.4, cuando, de hecho, p difiere de $1/2$, entonces n se agrandará más de lo necesario para el grado de confianza específico y, como resultado, se incrementará nuestro grado de confianza.

Teorema 9.5: Si utilizamos \hat{p} como un estimado de p , podemos tener, **al menos**, un $100(1 - \alpha)\%$ de confianza en que el error no excederá a una cantidad específica e cuando el tamaño de la muestra sea

$$n = \frac{z_{\alpha/2}^2}{4e^2}.$$

Ejemplo 9.16: ¿Qué tan grande debe ser una muestra en el ejemplo 9.14 si queremos tener al menos un 95% de confianza en que nuestra estimación de p está dentro de 0.02 del valor verdadero?

Solución: A diferencia del ejemplo 9.15, supondremos ahora que no se tomó una muestra preliminar para obtener una estimación de p . En consecuencia, podemos tener al menos un 95% de confianza en que nuestra proporción de la muestra no diferirá de la proporción verdadera en más de 0.02, si elegimos una muestra de tamaño

$$n = \frac{(1.96)^2}{(4)(0.02)^2} = 2401.$$

Si comparamos los resultados de los ejemplos 9.15 y 9.16, vemos que la información concerniente a p , proporcionada por una muestra preliminar, o quizás obtenida a partir de la experiencia, nos permite elegir una muestra más pequeña a la vez que mantenemos el grado de precisión requerido. ■

9.11 Dos muestras: estimación de la diferencia entre dos proporciones

Considere el problema en el que se busca estimar la diferencia entre dos parámetros binomiales p_1 y p_2 . Por ejemplo, p_1 podría ser la proporción de fumadores con cáncer de pulmón y p_2 la proporción de no fumadores con cáncer de pulmón, y el problema consistiría en estimar la diferencia entre estas dos proporciones. Primero seleccionamos muestras aleatorias independientes de tamaños n_1 y n_2 a partir de las dos poblaciones binomiales con medias $n_1 p_1$ y $n_2 p_2$, y varianzas $n_1 p_1 q_1$ y $n_2 p_2 q_2$, respectivamente, después determinamos los números x_1 y x_2 de personas con cáncer de pulmón en cada muestra, y formamos las proporciones $\hat{p}_1 = x_1/n_1$ y $\hat{p}_2 = x_2/n_2$. El estadístico $\hat{p}_1 - \hat{p}_2$ provee un estimador puntual de la diferencia entre las dos proporciones, $p_1 - p_2$. Por lo tanto, la diferencia de las proporciones muestrales, $\hat{p}_1 - \hat{p}_2$, se utilizará como la estimación puntual de $p_1 - p_2$.

Se puede establecer un intervalo de confianza para $p_1 - p_2$ considerando la distribución muestral de $\hat{p}_1 - \hat{p}_2$. De la sección 9.10 sabemos que \hat{p}_1 y \hat{p}_2 están distribuidos cada uno de forma aproximadamente normal, con medias p_1 y p_2 , y varianzas $p_1 q_1/n_1$ y $p_2 q_2/n_2$, respectivamente. Al elegir muestras independientes de las dos poblaciones nos aseguramos de que las variables \hat{p}_1 y \hat{p}_2 serán independientes y luego, por la propiedad reproductiva de la distribución normal que se estableció en el teorema 7.11, concluimos que $\hat{p}_1 - \hat{p}_2$ está distribuido de forma aproximadamente normal con media

$$\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$$

y varianza

$$\sigma_{\hat{p}_1 - \hat{p}_2}^2 = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}.$$

Por lo tanto, podemos asegurar que

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha,$$

donde

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{p_1 q_1/n_1 + p_2 q_2/n_2}}$$

y $z_{\alpha/2}$ es un valor por arriba del cual encontramos una área de $\alpha/2$ debajo de la curva normal estándar. Al sustituir para Z escribimos

$$P \left[-z_{\alpha/2} < \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{p_1 q_1/n_1 + p_2 q_2/n_2}} < z_{\alpha/2} \right] = 1 - \alpha.$$

Después de realizar las operaciones matemáticas usuales reemplazamos p_1, p_2, q_1 y q_2 bajo el signo de radical por sus estimaciones $\hat{p}_1 = x_1/n_1$, $\hat{p}_2 = x_2/n_2$, $\hat{q}_1 = 1 - \hat{p}_1$ y $\hat{q}_2 = 1 - \hat{p}_2$, siempre y cuando $n_1 \hat{p}_1$, $n_2 \hat{q}_1$, $n_2 \hat{p}_2$ y $n_1 \hat{q}_2$ sean todas mayores que o iguales a 5, y se obtiene el siguiente intervalo de confianza aproximado del $100(1 - \alpha)\%$ para $p_1 - p_2$:

Intervalo de confianza para $p_1 - p_2$ de una muestra grande

Si \hat{p}_1 y \hat{p}_2 son las proporciones de éxitos en muestras aleatorias de tamaños n_1 y n_2 , respectivamente, $\hat{q}_1 = 1 - \hat{p}_1$ y $\hat{q}_2 = 1 - \hat{p}_2$, un intervalo de confianza aproximado del $100(1 - \alpha)\%$ para la diferencia de dos parámetros binomiales $p_1 - p_2$ es dado por

$$(\hat{p}_1 - \hat{p}_2) - z_{\alpha/2} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}} < p_1 - p_2 < (\hat{p}_1 - \hat{p}_2) + z_{\alpha/2} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}},$$

donde $z_{\alpha/2}$ es el valor z que deja una área de $\alpha/2$ a la derecha.

Ejemplo 9.17: Se considera hacer un cierto cambio en el proceso de fabricación de partes componentes. Para determinar si el cambio en el proceso da como resultado una mejora, se toman muestras de partes fabricadas con el proceso nuevo y con el actual. Si se encuentra que 75 de 1500 artículos manufacturados con el proceso actual están defectuosos y 80 de 2000 manufacturados con el proceso nuevo también lo están, calcule un intervalo de confianza del 90% para la diferencia verdadera en la proporción de partes defectuosas entre el proceso actual y el nuevo.

Solución: Suponga que p_1 y p_2 son las proporciones verdaderas de partes defectuosas para los procesos actual y nuevo, respectivamente. En consecuencia, $\hat{p}_1 = 75/1500 = 0.05$ y $\hat{p}_2 = 80/2000 = 0.04$, y la estimación puntual de $p_1 - p_2$ es

$$\hat{p}_1 - \hat{p}_2 = 0.05 - 0.04 = 0.01.$$

Si utilizamos la tabla A.3, encontramos $z_{0.05} = 1.645$. Por lo tanto, al sustituir en la fórmula

$$1.645 \sqrt{\frac{(0.05)(0.95)}{1500} + \frac{(0.04)(0.96)}{2000}} = 0.0117,$$

encontramos que el intervalo de confianza del 90% es $-0.0017 < p_1 - p_2 < 0.0217$. Como el intervalo contiene el valor 0, no hay razón para creer que el nuevo proceso, comparado con el actual, disminuye en forma significativa la proporción de artículos defectuosos. ■

Hasta aquí todos los intervalos de confianza presentados son de la forma

estimación puntual $\pm K$ e.e.(estimación puntual),

donde K es una constante (ya sea t o el punto porcentual normal). Esta forma es válida cuando el parámetro es una media, una diferencia entre medias, una proporción o una diferencia entre proporciones, debido a la simetría de las distribuciones t y Z . Sin embargo, no se extiende a las varianzas ni a los cocientes de las varianzas, las cuales se examinarán en las secciones 9.12 y 9.13.

Ejercicios

En este conjunto de ejercicios, para una estimación respecto a una proporción, utilice sólo el método 1 para calcular los intervalos de confianza, a menos que se especifique otra cosa.

9.51 En una muestra aleatoria de 1000 viviendas en cierta ciudad se encuentra que 228 utilizan petróleo como combustible para la calefacción. Calcule intervalos de confianza del 99% para la proporción de viviendas en esta ciudad que utilizan petróleo con el fin mencionado. Utilice los dos métodos que se presentaron en la página 297.

9.52 Calcule intervalos de confianza del 95% para la proporción de artículos defectuosos que resultan de un proceso cuando se encuentra que una muestra de tamaño 100 produce 8 defectuosos. Utilice los dos métodos que se presentaron en la página 297.

9.53 a) Se selecciona una muestra aleatoria de 200 votantes en una ciudad y se encuentra que 114 apoyan un juicio de anexión. Calcule el intervalo de confianza del 96% para la parte de la población votante que está a favor del juicio.

b) ¿Qué podemos afirmar con 96% de confianza acerca de la posible magnitud de nuestro error, si estimamos que la fracción de votantes que está a favor del juicio de anexión es 0.57?

9.54 Un fabricante de reproductores de MP3 utiliza un conjunto de pruebas exhaustivas para evaluar el funcionamiento eléctrico de su producto. Todos los reproductores de MP3 deben pasar todas las pruebas antes de ser puestos a la venta. De una muestra aleatoria de 500 reproductores, 15 no pasan una o más de las pruebas. Calcule un intervalo de confianza del 90% para la proporción de los reproductores de MP3 de la población que pasan todas las pruebas.

9.55 Se está considerando un nuevo sistema de lanzamiento de cohetes para el despliegue de cohetes pequeños, de corto alcance. La probabilidad de que el sistema existente tenga un lanzamiento exitoso se representa con $p = 0.8$. Se toma una muestra de 40 lanzamientos experimentales con el nuevo sistema y 34 resultan exitosos.

- a) Construya un intervalo de confianza del 95% para p .
b) ¿Con base en sus resultados, concluiría que el nuevo sistema es mejor?

9.56 Un genetista está interesado en determinar la proporción de hombres africanos que padecen cierto trastorno sanguíneo menor. En una muestra aleatoria de 100 hombres africanos encuentra que 24 lo padecen.

- a) Calcule un intervalo de confianza del 99% para la proporción de hombres africanos que padecen este trastorno sanguíneo.

b) ¿Qué podríamos afirmar con 99% de confianza acerca de la posible magnitud de nuestro error, si estimamos que la proporción de hombres africanos con dicho trastorno sanguíneo es 0.24?

9.57 a) De acuerdo con un reporte del *Roanoke Times & World-News*, aproximadamente $2/3$ de los 1600 adultos encuestados vía telefónica dijeron que piensan que invertir en el programa del transbordador espacial es bueno para Estados Unidos. Calcule un intervalo de confianza del 95% para la proporción de adultos estadounidenses que piensan que el programa del transbordador espacial es una buena inversión para su país.

b) ¿Qué podríamos afirmar con un 95% de confianza acerca de la posible magnitud de nuestro error, si estimamos que la proporción de adultos estadounidenses que piensan que el programa del transbordador espacial es una buena inversión es de $2/3$?

9.58 En el artículo del periódico al que se hace referencia en el ejercicio 9.57, 32% de los 1600 adultos encuestados dijo que el programa espacial estadounidense debería enfatizar la exploración científica. ¿Qué tamaño debería tener una muestra de adultos para la encuesta si se desea tener un 95% de confianza en que el porcentaje estimado esté dentro del 2% del porcentaje verdadero?

9.59 ¿Qué tamaño debería tener una muestra si deseamos tener un 96% de confianza en que nuestra proporción de la muestra en el ejercicio 9.53 esté dentro del 0.02 de la fracción verdadera de la población votante?

9.60 ¿Qué tamaño debería tener una muestra si deseamos tener un 99% de confianza en que nuestra proporción de la muestra en el ejercicio 9.51 esté dentro del 0.05 de la proporción verdadera de viviendas en esa ciudad que utilizan petróleo como combustible para la calefacción?

9.61 ¿Qué tamaño debería tener una muestra en el ejercicio 9.52 si deseamos tener un 98% de confianza en que nuestra proporción de la muestra esté dentro del 0.05 de la proporción verdadera de defectuosos?

9.62 Una conjetura de un catedrático del departamento de microbiología, de la Facultad de Odontología de la Universidad de Washington, en St. Louis, Missouri, afirma que un par de tasas diarias de té verde o negro proporciona suficiente flúor para evitar el deterioro de los dientes. ¿Qué tan grande debería ser la muestra para estimar el porcentaje de habitantes de cierta ciudad que están a favor de tener agua fluorada, si se desea tener al menos un 99% de confianza en que el estimado está dentro del 1% del porcentaje verdadero?

9.63 Se llevará a cabo un estudio para estimar el porcentaje de ciudadanos de una ciudad que están a favor de tener agua fluorada. ¿Qué tan grande debería ser la muestra si se desea tener al menos 95% de confianza en que el estimado esté dentro del 1% del porcentaje verdadero?

9.64 Se realizará un estudio para estimar la proporción de residentes de cierta ciudad y sus suburbios que está a favor de que se construya una planta de energía nuclear cerca de la ciudad. ¿Qué tan grande debería ser la muestra, si se desea tener al menos un 95% de confianza en que el estimado esté dentro del 0.04 de la verdadera proporción de residentes que están a favor de que se construya la planta de energía nuclear?

9.65 A cierto genetista le interesa determinar la proporción de hombres y mujeres de la población que padecen cierto trastorno sanguíneo menor. En una muestra aleatoria de 1000 hombres encuentra que 250 lo padecen; mientras que de 1000 mujeres examinadas, 275 parecen padecerlo. Calcule un intervalo de confianza del 95% para la diferencia entre la proporción de hombres y mujeres que padecen el trastorno sanguíneo.

9.66 Se encuestan 10 escuelas de ingeniería de Estados Unidos. La muestra contiene a 250 ingenieros eléctricos, de los cuales 80 son mujeres; y 175 ingenieros químicos, de los cuales 40 son mujeres. Calcule un intervalo de confianza del 90% para la diferencia entre la proporción de mujeres en estos dos campos de la ingeniería. ¿Hay una diferencia significativa entre las dos proporciones?

9.67 Se llevó a cabo una prueba clínica para determinar si cierto tipo de vacuna tiene un efecto sobre la incidencia de cierta enfermedad. Una muestra de 1000 ratas, 500 de las cuales recibieron la vacuna, se mantuvo en un ambiente controlado durante un periodo de un

año. En el grupo que no fue vacunado, 120 ratas presentaron la enfermedad, mientras que en el grupo inoculado 98 ratas la contrajeron. Si p_1 es la probabilidad de incidencia de la enfermedad en las ratas sin vacuna y p_2 es la probabilidad de incidencia en las ratas inoculadas, calcule un intervalo de confianza del 90% para $p_1 - p_2$.

9.68 En el estudio *Germination and Emergence of Broccoli*, realizado por el Departamento de horticultura de la Virginia Tech, un investigador encontró que a 5°C, de 20 semillas de brócoli germinaron 10; en tanto que a 15°C, de 20 semillas germinaron 15. Calcule un intervalo de confianza del 95% para la diferencia en la proporción de semillas que germinaron a las dos temperaturas y decida si esta diferencia es significativa.

9.69 Una encuesta de 1000 estudiantes reveló que 274 eligen al equipo profesional de beisbol A como su equipo favorito. En 1991 se realizó una encuesta similar con 760 estudiantes y 240 de ellos también eligieron a ese equipo como su favorito. Calcule un intervalo de confianza del 95% para la diferencia entre la proporción de estudiantes que favorecen al equipo A en las dos encuestas. ¿Hay una diferencia significativa?

9.70 De acuerdo con el *USA Today* (17 de marzo de 1997), las mujeres constituían el 33.7% del personal de redacción en las estaciones locales de televisión en 1990 y el 36.2% en 1994. Suponga que en 1990 y en 1994 se contrataron 20 nuevos empleados para el personal de redacción.

- Estime el número de trabajadores que habrían sido mujeres en 1990 y en 1994, respectivamente.
- Calcule un intervalo de confianza del 95% para saber si hay evidencia de que la proporción de mujeres contratadas para el equipo de redacción fue mayor en 1994 que en 1990.

9.12 Una sola muestra: estimación de la varianza

Si extraemos una muestra de tamaño n de una población normal con varianza σ^2 y calculamos la varianza muestral s^2 , obtenemos un valor del estadístico S^2 . Esta varianza muestral calculada se utiliza como una estimación puntual de σ^2 . En consecuencia, al estadístico S^2 se le denomina estimador de σ^2 .

Se puede establecer una estimación por intervalos de σ^2 utilizando el estadístico

$$X^2 = \frac{(n-1)S^2}{\sigma^2}.$$

De acuerdo con el teorema 8.4, cuando las muestras se toman de una población normal el estadístico X^2 tiene una distribución chi cuadrada con $n-1$ grados de libertad. Podemos escribir (véase la figura 9.7)

$$P(\chi_{1-\alpha/2}^2 < X^2 < \chi_{\alpha/2}^2) = 1 - \alpha.$$

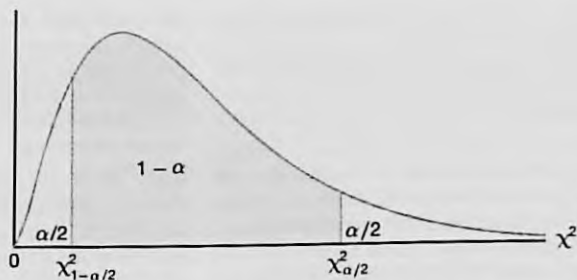


Figura 9.7: $P(\chi^2_{1-\alpha/2} < X^2 < \chi^2_{\alpha/2}) = 1 - \alpha$

donde $\chi^2_{1-\alpha/2}$ y $\chi^2_{\alpha/2}$ son valores de la distribución chi cuadrada con $n - 1$ grados de libertad, que dejan áreas de $1 - \alpha/2$ y $\alpha/2$, respectivamente, a la derecha. Al sustituir para X^2 escribimos

$$P \left[\chi^2_{1-\alpha/2} < \frac{(n-1)S^2}{\sigma^2} < \chi^2_{\alpha/2} \right] = 1 - \alpha.$$

Si dividimos cada término de la desigualdad entre $(n-1)S^2$, y después invertimos cada término (lo que cambia el sentido de las desigualdades), obtenemos

$$P \left[\frac{(n-1)S^2}{\chi^2_{\alpha/2}} < \sigma^2 < \frac{(n-1)S^2}{\chi^2_{1-\alpha/2}} \right] = 1 - \alpha.$$

Para una muestra aleatoria de tamaño n , tomada de una población normal, se calcula la varianza muestral s^2 y se obtiene el siguiente intervalo de confianza del $100(1 - \alpha)\%$ para σ^2 .

Intervalo de confianza para σ^2 Si s^2 es la varianza de una muestra aleatoria de tamaño n de una población normal, un intervalo de confianza del $100(1 - \alpha)\%$ para σ^2 es

$$\frac{(n-1)s^2}{\chi^2_{\alpha/2}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{1-\alpha/2}},$$

donde $\chi^2_{\alpha/2}$ y $\chi^2_{1-\alpha/2}$ son valores χ^2 con $v = n - 1$ grados de libertad, que dejan áreas de $\alpha/2$ y $1 - \alpha/2$, respectivamente, a la derecha.

Un intervalo de confianza aproximado a $100(1 - \alpha)\%$ para σ se obtiene tomando la raíz cuadrada de cada extremo del intervalo para σ^2 .

Ejemplo 9.18: Los siguientes son los pesos, en decagramos, de 10 paquetes de semillas de pasto distribuidas por cierta empresa: 46.4, 46.1, 45.8, 47.0, 46.1, 45.9, 45.8, 46.9, 45.2 y 46.0. Calcule un intervalo de confianza del 95% para la varianza de todos los pesos de este tipo de paquetes de semillas de pasto distribuidos por la empresa. Suponga una población normal.

Solución: Primero calculamos

$$\begin{aligned} s^2 &= \frac{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}{n(n-1)} \\ &= \frac{(10)(21,273.12) - (461.2)^2}{(10)(9)} = 0.286. \end{aligned}$$

Para obtener un intervalo de confianza del 95% elegimos $\alpha = 0.05$. Después, usando la tabla A.5 con $\nu = 9$ grados de libertad, encontramos $\chi_{0.025}^2 = 19.023$ y $\chi_{0.975}^2 = 2.700$. Por lo tanto, el intervalo de confianza del 95% para σ^2 es

$$\frac{(9)(0.286)}{19.023} < \sigma^2 < \frac{(9)(0.286)}{2.700},$$

o simplemente $0.135 < \sigma^2 < 0.953$. ▮

9.13 Dos muestras: estimación de la proporción de dos varianzas

Una estimación puntual de la proporción de dos varianzas de la población σ_1^2/σ_2^2 es dada por la proporción s_1^2/s_2^2 de las varianzas muestrales. En consecuencia, el estadístico S_1^2/S_2^2 se conoce como un estimador de σ_1^2/σ_2^2 .

Si σ_1^2 y σ_2^2 son las varianzas de poblaciones normales, podemos establecer una estimación por intervalos de σ_1^2/σ_2^2 usando el estadístico

$$F = \frac{\sigma_2^2 S_1^2}{\sigma_1^2 S_2^2}.$$

De acuerdo con el teorema 8.8, la variable aleatoria F tiene una distribución F con $\nu_1 = n_1 - 1$ y $\nu_2 = n_2 - 1$ grados de libertad. Por lo tanto, podemos escribir (véase la figura 9.8)

$$P[f_{1-\alpha/2}(\nu_1, \nu_2) < F < f_{\alpha/2}(\nu_1, \nu_2)] = 1 - \alpha,$$

donde $f_{1-\alpha/2}(\nu_1, \nu_2)$ y $f_{\alpha/2}(\nu_1, \nu_2)$ son los valores de la distribución F con ν_1 y ν_2 grados de libertad, que dejan áreas de $1 - \alpha/2$ y $\alpha/2$, respectivamente, a la derecha.

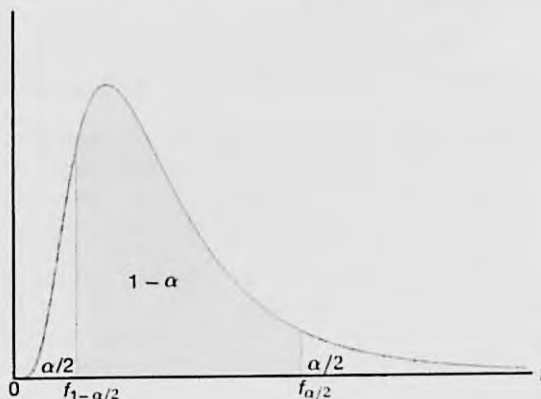


Figura 9.8: $P[f_{1-\alpha/2}(\nu_1, \nu_2) < F < f_{\alpha/2}(\nu_1, \nu_2)] = 1 - \alpha$.

Al sustituir para F , escribimos

$$P \left[f_{1-\alpha/2}(v_1, v_2) < \frac{\sigma_2^2 S_1^2}{\sigma_1^2 S_2^2} < f_{\alpha/2}(v_1, v_2) \right] = 1 - \alpha.$$

Si multiplicamos cada término de la desigualdad por S_2^2/S_1^2 , y después invertimos cada término, obtenemos

$$P \left[\frac{S_1^2}{S_2^2} \frac{1}{f_{\alpha/2}(v_1, v_2)} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{S_1^2}{S_2^2} \frac{1}{f_{1-\alpha/2}(v_1, v_2)} \right] = 1 - \alpha.$$

Los resultados del teorema 8.7 nos permiten reemplazar la cantidad $f_{1-\alpha/2}(v_1, v_2)$ por $1/f_{\alpha/2}(v_2, v_1)$. Por lo tanto,

$$P \left[\frac{S_1^2}{S_2^2} \frac{1}{f_{\alpha/2}(v_1, v_2)} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{S_1^2}{S_2^2} f_{\alpha/2}(v_2, v_1) \right] = 1 - \alpha.$$

Para cualesquiera dos muestras aleatorias independientes de tamaño n_1 y n_2 que se seleccionan de dos poblaciones normales, se calcula la proporción de las varianzas muestrales s_1^2/s_2^2 y se obtiene el siguiente intervalo de confianza del $100(1 - \alpha)\%$ para σ_1^2/σ_2^2 .

Intervalo de confianza para σ_1^2/σ_2^2 Si s_1^2 y s_2^2 son las varianzas de muestras independientes de tamaño n_1 y n_2 , respectivamente, tomadas de poblaciones normales, entonces un intervalo de confianza del $100(1 - \alpha)\%$ para σ_1^2/σ_2^2 es

$$\frac{s_1^2}{s_2^2} \frac{1}{f_{\alpha/2}(v_1, v_2)} < \frac{\sigma_1^2}{\sigma_2^2} < \frac{s_1^2}{s_2^2} f_{\alpha/2}(v_2, v_1),$$

donde $f_{\alpha/2}(v_1, v_2)$ es un valor f con $v_1 = n_1 - 1$ y $v_2 = n_2 - 1$ grados de libertad que deja una área de $\alpha/2$ a la derecha, y $f_{\alpha/2}(v_2, v_1)$ es un valor f similar con $v_2 = n_2 - 1$ y $v_1 = n_1 - 1$ grados de libertad.

Como vimos en la sección 9.12, tomando la raíz cuadrada de cada extremo del intervalo para σ_1^2/σ_2^2 , se obtiene un intervalo de confianza del $100(1 - \alpha)\%$ para σ_1/σ_2 .

Ejemplo 9.19: En el ejemplo 9.12 de la página 290 se construyó un intervalo de confianza para la diferencia en el contenido medio de ortofósforo de dos estaciones ubicadas sobre el río James, medido en miligramos por litro, suponiendo que las varianzas normales de la población son diferentes. Justifique esta suposición construyendo intervalos de confianza del 98% para σ_1^2/σ_2^2 y para σ_1/σ_2 , donde σ_1^2 y σ_2^2 son las varianzas de la población del contenido de ortofósforo en la estación 1 y en la estación 2, respectivamente.

Solución: Del ejemplo 9.12 tenemos $n_1 = 15$, $n_2 = 12$, $s_1 = 3.07$ y $s_2 = 0.80$. Para un intervalo de confianza del 98%, $\alpha = 0.02$. Al interpolar en la tabla A.6 encontramos $f_{0.01}(14, 11) \approx 4.30$ y $f_{0.01}(11, 14) \approx 3.87$. Por lo tanto, el intervalo de confianza del 98% para σ_1^2/σ_2^2 es

$$\left(\frac{3.07^2}{0.80^2} \right) \left(\frac{1}{4.30} \right) < \frac{\sigma_1^2}{\sigma_2^2} < \left(\frac{3.07^2}{0.80^2} \right) (3.87),$$

que se simplifica a $3.425 < \frac{\sigma_1^2}{\sigma_2^2} < 56.991$. Al calcular las raíces cuadradas de los límites de confianza encontramos que un intervalo de confianza del 98% para σ_1/σ_2 es

$$1.851 < \frac{\sigma_1}{\sigma_2} < 7.549.$$

Como este intervalo no permite la posibilidad de que σ_1/σ_2 sea igual a 1, es correcto suponer que $\sigma_1 \neq \sigma_2$ o $\sigma_1^2 \neq \sigma_2^2$ en el ejemplo 9.12. ■

Ejercicios

9.71 Un fabricante de baterías para automóvil afirma que sus baterías durarán, en promedio, 3 años con una varianza de 1 año. Suponga que 5 de estas baterías tienen duraciones de 1.9, 2.4, 3.0, 3.5 y 4.2 años y con base en esto construya un intervalo de confianza del 95% para σ^2 , después decida si la afirmación del fabricante de que $\sigma^2 = 1$ es válida. Suponga que la población de duraciones de las baterías se distribuye de forma aproximadamente normal.

9.72 Una muestra aleatoria de 20 estudiantes obtuvo una media de $\bar{x} = 72$ y una varianza de $s^2 = 16$ en un examen universitario de colocación en matemáticas. Suponga que las calificaciones se distribuyen normalmente y con base en esto construya un intervalo de confianza del 98% para σ^2 .

9.73 Construya un intervalo de confianza del 95% para σ^2 en el ejercicio 9.9 de la página 283.

9.74 Construya un intervalo de confianza del 99% para σ^2 en el ejercicio 9.11 de la página 283.

9.75 Construya un intervalo de confianza del 99% para σ en el ejercicio 9.12 de la página 283.

9.76 Construya un intervalo de confianza del 90% para σ en el ejercicio 9.13 de la página 283.

9.77 Construya un intervalo de confianza del 98% para σ_1/σ_2 en el ejercicio 9.42 de la página 295, donde σ_1 y σ_2 son, respectivamente, las desviaciones estándar para las distancias recorridas por litro de combustible de los camiones compactos Volkswagen y Toyota.

9.78 Construya un intervalo de confianza del 90% para σ_1^2/σ_2^2 en el ejercicio 9.43 de la página 295. ¿Se justifica que supongamos que $\sigma_1^2 \neq \sigma_2^2$ cuando construimos nuestro intervalo de confianza para $\mu_1 - \mu_2$?

9.79 Construya un intervalo de confianza del 90% para σ_1^2/σ_2^2 en el ejercicio 9.46 de la página 295. ¿Deberíamos suponer que $\sigma_1^2 = \sigma_2^2$ cuando construimos nuestro intervalo de confianza para $\mu_1 - \mu_2$?

9.80 Construya un intervalo de confianza del 95% para σ_A^2/σ_B^2 en el ejercicio 9.49 de la página 295. ¿Tendría que utilizar la suposición de la igualdad de la varianza?

9.14 Estimación de la máxima verosimilitud (opcional)

A menudo los estimadores de parámetros han tenido que recurrir a la intuición. El estimador \bar{X} ciertamente parece razonable como estimador de una media de la población μ . La virtud de S^2 como estimador de σ^2 se destaca en el estudio de estimadores insesgados de la sección 9.3. El estimador para un parámetro binomial p es simplemente una proporción de la muestra que, desde luego, es un *promedio* y recurre al sentido común. Sin embargo, hay muchas situaciones en las que no es del todo evidente cuál debería ser el estimador adecuado. Como resultado, el estudiante de estadística tiene mucho que aprender respecto a las diferentes filosofías que producen distintos métodos de estimación. En esta sección estudiaremos el **método de máxima verosimilitud**.

La estimación por máxima verosimilitud representa uno de los métodos de estimación más importantes en toda la estadística inferencial. No explicaremos el método de manera detallada; más bien, intentaremos transmitir la filosofía de la máxima verosimilitud e ilustrarla con ejemplos que la relacionan con otros problemas de estimación que se examinan en este capítulo.

Función de verosimilitud

Como el nombre lo indica, el método de máxima verosimilitud es aquel para el que se maximiza la *función de verosimilitud*, lo cual se ilustra mejor con un ejemplo que incluye una distribución discreta y un solo parámetro. Consideremos que X_1, X_2, \dots, X_n son las variables aleatorias independientes tomadas de una distribución de probabilidad discreta representada por $f(x, \theta)$, donde θ es un solo parámetro de la distribución. Ahora bien,

$$\begin{aligned} L(x_1, x_2, \dots, x_n; \theta) &= f(x_1, x_2, \dots, x_n; \theta) \\ &= f(x_1, \theta) f(x_2, \theta) \cdots f(x_n, \theta) \end{aligned}$$

es la *distribución conjunta de las variables aleatorias*, la cual a menudo se denomina función de probabilidad. Observe que la variable de la función de probabilidad es θ , no x . Represente con x_1, x_2, \dots, x_n los valores observados en una muestra. En el caso de una variable aleatoria discreta, la interpretación es muy clara. La cantidad $L(x_1, x_2, \dots, x_n; \theta)$, la *verosimilitud de la muestra*, es la siguiente probabilidad conjunta:

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | \theta),$$

que es la probabilidad de obtener los valores muestrales x_1, x_2, \dots, x_n . Para el caso discreto el estimador de máxima verosimilitud es el que da como resultado un valor máximo para esta probabilidad conjunta, o el que maximiza la probabilidad de la muestra.

Considere un ejemplo ficticio en el cual se inspeccionan tres artículos que salen de una línea de ensamble. Los artículos se clasifican como defectuosos o no defectuosos, de manera que se aplica el proceso de Bernoulli. La inspección de los tres artículos da como resultado dos artículos no defectuosos seguidos por uno defectuoso. Nos interesa estimar p , la proporción de artículos no defectuosos en el proceso. La probabilidad de la muestra para este ejemplo es dada por

$$p \cdot p \cdot q = p^2 q = p^2 - p^3,$$

donde $q = 1 - p$. La estimación de máxima verosimilitud daría un estimado de p para el que se maximiza la verosimilitud. Resulta claro que si diferenciamos la verosimilitud respecto a p , igualamos la derivada a cero y la resolvemos, obtenemos el valor

$$\hat{p} = \frac{2}{3}.$$

Entonces, desde luego, en esta situación $\hat{p} = 2/3$ es la proporción muestral defectuosa y, por ello, un estimador razonable de la probabilidad de un artículo defectuoso. El lector debería intentar comprender que la filosofía de la estimación de máxima verosimilitud proviene de la noción de que el estimador razonable de un parámetro que se basa en información muestral *es el valor del parámetro que produce la mayor probabilidad de obtener la muestra*. Ésta es, de hecho, la interpretación para el caso discreto, ya que la verosimilitud es la probabilidad de observar de manera conjunta los valores en la muestra.

Así, mientras que la interpretación de la función de verosimilitud como una probabilidad conjunta se limita al caso discreto, la noción de máxima verosimilitud se extiende a la estimación de parámetros de una distribución continua. Presentamos ahora una definición formal de la estimación de máxima verosimilitud.

Definición 9.3: Dadas las observaciones independientes x_1, x_2, \dots, x_n de una función de densidad de probabilidad (caso continuo) o de una función de masa de probabilidad (caso discreto) $f(x, \theta)$, el estimador de máxima verosimilitud $\hat{\theta}$ es el que maximiza la función de probabilidad

$$L(x_1, x_2, \dots, x_n; \theta) = f(x; \theta) = f(x_1, \theta) f(x_2, \theta) \cdots f(x_n, \theta).$$

Muy a menudo conviene trabajar con el logaritmo natural de la función de verosimilitud para encontrar el máximo de esa función. Considere el siguiente ejemplo acerca del parámetro μ de una distribución de Poisson.

Ejemplo 9.20: Considere una distribución de Poisson con la siguiente función de masa de probabilidad

$$f(x|\mu) = \frac{e^{-\mu} \mu^x}{x!}, \quad x = 0, 1, 2, \dots$$

Suponga que se toma una muestra aleatoria x_1, x_2, \dots, x_n de la distribución. ¿Cuál es la estimación de máxima verosimilitud de μ ?

Solución: La función de probabilidad es

$$L(x_1, x_2, \dots, x_n; \mu) = \prod_{i=1}^n f(x_i|\mu) = \frac{e^{-n\mu} \mu^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}.$$

Considere ahora

$$\begin{aligned} \ln L(x_1, x_2, \dots, x_n; \mu) &= -n\mu + \sum_{i=1}^n x_i \ln \mu - \ln \prod_{i=1}^n x_i! \\ \frac{\partial \ln L(x_1, x_2, \dots, x_n; \mu)}{\partial \mu} &= -n + \sum_{i=1}^n \frac{x_i}{\mu}. \end{aligned}$$

Resolver para $\hat{\mu}$, el estimador de máxima verosimilitud, implica definir la derivada para cero y resolver para el parámetro. Por consiguiente,

$$\hat{\mu} = \sum_{i=1}^n \frac{x_i}{n} = \bar{x}.$$

La segunda derivada de la función de verosimilitud logarítmica es negativa, lo cual implica que la solución anterior realmente es un máximo. Como μ es la media de la distribución de Poisson (capítulo 5), el promedio muestral en realidad parecería ser un estimador razonable. \blacksquare

El siguiente ejemplo presenta el uso del método de máxima verosimilitud para calcular estimados de dos parámetros. Simplemente encontramos los valores de los parámetros que maximizan (de forma conjunta) la función de probabilidad.

Ejemplo 9.21: Considere una muestra aleatoria x_1, x_2, \dots, x_n de una distribución normal $N(\mu, \sigma)$. Calcule los estimadores de máxima verosimilitud para μ y σ^2 .

Solución: La función de verosimilitud para la distribución normal es

$$L(x_1, x_2, \dots, x_n; \mu, \sigma^2) = \frac{1}{(2\pi)^{n/2} (\sigma^2)^{n/2}} \exp \left[-\frac{1}{2} \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma} \right)^2 \right].$$

Al usar logaritmos obtenemos

$$\ln L(x_1, x_2, \dots, x_n; \mu, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2} \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma} \right)^2.$$

Por lo tanto,

$$\frac{\partial \ln L}{\partial \mu} = \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma^2} \right)$$

y

$$\frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2.$$

Al igualar ambas derivadas a cero, obtenemos

$$\sum_{i=1}^n x_i - n\mu = 0 \quad \text{y} \quad n\sigma^2 = \sum_{i=1}^n (x_i - \mu)^2.$$

Por consiguiente, el estimador de máxima verosimilitud de μ es dado por

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x},$$

que es un resultado satisfactorio, ya que \bar{x} ha desempeñado un papel tan importante en este capítulo como un estimador puntual de μ . Por otro lado, el estimador de máxima verosimilitud de σ^2 es

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Al verificar la matriz derivada parcial de segundo orden se confirma que la solución da como resultado el máximo de la función de verosimilitud. \blacksquare

Resulta interesante notar la distinción entre el estimador de máxima verosimilitud de σ^2 y el estimador insesgado S^2 que se presentó al principio de este capítulo. Los numeradores son idénticos, desde luego, y el denominador lo constituyen los "grados de libertad" $n - 1$ para el estimador insesgado, y n para el estimador de máxima verosimilitud. Los estimadores de máxima verosimilitud no necesariamente gozan de la propiedad de carecer de sesgo. Sin embargo, los estimadores de máxima verosimilitud tienen importantes propiedades asintóticas.

Ejemplo 9.22: Suponga que en un estudio biomédico se utilizan 10 ratas a las que después de inyectarles células cancerosas se les suministra un fármaco contra el cáncer diseñado para aumentar su tasa de supervivencia. Los tiempos de supervivencia, en meses, son 14, 17, 27, 18, 12,

8, 22, 13, 19 y 12. Suponga que se trata de una distribución exponencial. Calcule un estimado de máxima verosimilitud de la supervivencia media.

Solución: Del capítulo 6 sabemos que la función de densidad de probabilidad para la variable aleatoria exponencial X es

$$f(x, \beta) = \begin{cases} \frac{1}{\beta} e^{-x/\beta}, & x > 0, \\ 0, & \text{en cualquier caso.} \end{cases}$$

Por consiguiente, la función de verosimilitud logarítmica de los datos, dado que $n = 10$, es

$$\ln L(x_1, x_2, \dots, x_{10}; \beta) = -10 \ln \beta - \frac{1}{\beta} \sum_{i=1}^{10} x_i.$$

Si se establece que

$$\frac{\partial \ln L}{\partial \beta} = -\frac{10}{\beta} + \frac{1}{\beta^2} \sum_{i=1}^{10} x_i = 0$$

implica que

$$\hat{\beta} = \frac{1}{10} \sum_{i=1}^{10} x_i = \bar{x} = 16.2.$$

Si se evalúa la segunda derivada de la función de verosimilitud logarítmica en el valor $\hat{\beta}$ anterior se produce un valor negativo. Como resultado, el estimador del parámetro β , la media de la población, es el promedio muestral \bar{x} . ▮

El siguiente ejemplo ilustra el estimador de máxima verosimilitud para una distribución que no se incluye en los capítulos anteriores.

Ejemplo 9.23: Se sabe que una muestra que consta de los valores 12, 11.2, 13.5, 12.3, 13.8 y 11.9 proviene de una población con la siguiente función de densidad

$$f(x; \theta) = \begin{cases} \frac{\theta}{x^{\theta+1}}, & x > 1, \\ 0, & \text{en cualquier caso,} \end{cases}$$

donde $\theta > 0$. Calcule la estimación de máxima verosimilitud de θ .

Solución: La función de verosimilitud de n observaciones de esta población se escribe como

$$L(x_1, x_2, \dots, x_{10}; \theta) = \prod_{i=1}^n \frac{\theta}{x_i^{\theta+1}} = \frac{\theta^n}{(\prod_{i=1}^n x_i)^{\theta+1}},$$

lo cual implica que

$$\ln L(x_1, x_2, \dots, x_{10}; \theta) = n \ln(\theta) - (\theta + 1) \sum_{i=1}^n \ln(x_i).$$

Si establecemos que $0 = \frac{\partial \ln L}{\partial \theta} = \frac{n}{\theta} - \sum_{i=1}^n \ln(x_i)$ da como resultado

$$\hat{\theta} = \frac{n}{\sum_{i=1}^n \ln(x_i)} = \frac{6}{\ln(12) + \ln(11.2) + \ln(13.5) + \ln(12.3) + \ln(13.8) + \ln(11.9)} = 0.3970.$$

Como la segunda derivada de L es $-n/\theta^2$, que siempre es negativa, la función de probabilidad alcanza su valor máximo en $\hat{\theta}$. J

Comentarios adicionales respecto a la estimación de máxima verosimilitud

Un análisis detallado de las propiedades de la estimación de máxima verosimilitud está fuera del alcance de este libro y, por lo general, es un tema importante en un curso teórico de estadística inferencial. El método de máxima verosimilitud permite al analista utilizar el conocimiento de la distribución para determinar un estimador adecuado. *El método de máxima verosimilitud no se puede aplicar si no se conoce la distribución subyacente.* En el ejemplo 9.21 aprendimos que el estimador de máxima verosimilitud no necesariamente carece de sesgo. El estimador de máxima verosimilitud es insesgado *asintóticamente* o *en el límite*; es decir, la magnitud del sesgo se aproxima a cero a medida que la muestra se hace más grande. Al principio de este capítulo examinamos la noción de eficacia, que se vincula con la propiedad de la varianza de un estimador. Los estimadores de máxima verosimilitud tienen propiedades de varianza deseables en el límite. El lector debería consultar la obra de Lehmann y D'Abbrera (1998) para más detalles.

Ejercicios

9.81 Suponga que hay n ensayos x_1, x_2, \dots, x_n de un proceso de Bernoulli con parámetro p , la probabilidad de un éxito. Esto es, la probabilidad de r éxitos es dada por $\binom{n}{r} p^r (1-p)^{n-r}$. Determine el estimador de máxima verosimilitud para el parámetro p .

9.82 Considere la distribución logarítmica normal con la función de densidad dada en la sección 6.9. Suponga que tiene una muestra aleatoria x_1, x_2, \dots, x_n de una distribución logarítmica normal.

- Escriba la función de verosimilitud.
- Desarrolle los estimadores de máxima verosimilitud de μ y σ^2 .

9.83 Considere una muestra aleatoria de x_1, \dots, x_n obtenida de la distribución gamma descrita en la sección 6.6. Suponga que conoce el parámetro α , el cual digamos que es 5, y con base en esto determine la estimación de máxima verosimilitud para el parámetro β .

9.84 Considere una muestra aleatoria de x_1, x_2, \dots, x_n observaciones de una distribución de Weibull con parámetros α y β , y la siguiente función de densidad

$$f(x) = \begin{cases} \alpha \beta x^{\beta-1} e^{-\alpha x^\beta}, & x > 0, \\ 0, & \text{en cualquier caso.} \end{cases}$$

para $\alpha, \beta > 0$.

- Escriba la función de verosimilitud.
- Escriba las ecuaciones que al resolverse proporcionan los estimadores de máxima verosimilitud de α y β .

9.85 Considere una muestra aleatoria de x_1, \dots, x_n obtenida de una distribución uniforme $U(0, \theta)$, con el parámetro θ desconocido, donde $\theta > 0$. Determine el estimador de máxima verosimilitud de θ .

9.86 Considere las observaciones independientes de x_1, x_2, \dots, x_n de la distribución gamma que se analizó en la sección 6.6.

- a) Escriba la función de verosimilitud.
 b) Escriba un conjunto de ecuaciones que, cuando se resuelven, proporcionan los estimadores de máxima verosimilitud de α y β .

9.87 Considere un experimento hipotético en el que un hombre que tiene un hongo utiliza un medicamento fungicida y se cura. Por lo tanto, considere que se trata de una muestra de una distribución de Bernoulli con la siguiente función de probabilidad

$$f(x) = p^x q^{1-x}, \quad x = 0, 1,$$

donde p es la probabilidad de un éxito (curación) y $q = 1 - p$. Ahora, desde luego, la información muestral da $x = 1$. Escriba un procedimiento que demuestre que $\hat{p} = 1.0$ es el estimador de máxima probabilidad de curación.

9.88 Considere la observación X de la distribución binomial negativa dada en la sección 5.4. Calcule el estimador de máxima verosimilitud para p , suponiendo que se conoce k .

Ejercicios de repaso

9.89 Considere dos estimadores de σ^2 para una muestra x_1, x_2, \dots, x_n que se extrae de una distribución normal con media μ y varianza σ^2 . Los estimadores son el estimador insesgado $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ y el estimador de máxima verosimilitud $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$.

Analice las propiedades de la varianza de estos dos estimadores.

9.90 De acuerdo con el *Roanoke Times*, McDonald's vendió 42.1% de la participación del mercado de hamburguesas. Una muestra aleatoria de 75 hamburguesas vendidas reveló que 28 de ellas fueron vendidas por McDonald's. Utilice el material de la sección 9.10 para determinar si esta información respalda la afirmación del *Roanoke Times*.

9.91 Se afirma que un individuo podrá reducir, en un lapso de 2 semanas, un promedio de 4.5 kilogramos de peso con una nueva dieta. Los pesos de 7 mujeres que siguieron esta dieta se registraron antes y después de un periodo de 2 semanas.

Mujer	Peso antes	Peso después
1	58.5	60.0
2	60.3	54.9
3	61.7	58.1
4	69.0	62.1
5	64.0	58.5
6	62.6	59.9
7	56.7	54.4

vera. El calcio es un elemento necesario para las plantas y los animales. La cantidad que la planta toma y almacena está estrechamente correlacionada con la cantidad presente en el suelo. Se formuló la hipótesis de que el fuego podría cambiar los niveles de calcio presentes en el suelo y, por lo tanto, influir en la cantidad disponible para los venados. Se seleccionó una extensión grande de tierra en el bosque Fishburn para provocar un incendio controlado. Justo antes de la quema se tomaron muestras de suelo de 12 parcelas con la misma área y se analizaron para verificar su contenido de calcio. Después del incendio se volvieron a analizar los niveles de calcio en las mismas parcelas. Los valores obtenidos, en kilogramos por parcela, se presentan en la siguiente tabla:

Parcela	Nivel de calcio (kg/parcela)	
	Antes del incendio	Después del incendio
1	50	9
2	50	18
3	82	45
4	64	18
5	82	18
6	73	9
7	77	32
8	54	9
9	23	18
10	45	9
11	36	9
12	54	9

Construya un intervalo de confianza del 95% para la diferencia media en los niveles de calcio presentes en el suelo antes y después del incendio controlado. Suponga que la distribución de las diferencias en los niveles de calcio es aproximadamente normal.

9.93 El dueño de un gimnasio afirma que una persona podrá reducir, en un periodo de 5 días, un promedio de 2 centímetros en su talla de cintura con un nuevo programa de ejercicios. En la siguiente tabla se presentan

Pruebe la afirmación sobre la dieta calculando un intervalo de confianza del 95% para la diferencia media en el peso. Suponga que las diferencias de los pesos se distribuyen de forma aproximadamente normal.

9.92 En Virginia Tech se realizó un estudio para determinar si se puede utilizar el fuego como una herramienta de control viable para aumentar la cantidad de forraje disponible para los venados durante los meses críticos a finales del invierno y principios de la prima-

las tallas de cintura de 6 hombres que participaron en este programa de ejercicios antes y después del periodo de 5 días:

Hombre	Talla de cintura antes	Talla de cintura después
1	90.4	91.7
2	95.5	93.9
3	98.7	97.4
4	115.9	112.8
5	104.0	101.3
6	85.6	84.0

Mediante el cálculo de un intervalo de confianza del 95% para la reducción media en la talla de cintura determine si la afirmación del dueño del gimnasio es válida. Suponga que la distribución de las diferencias en las tallas de cintura antes y después del programa es aproximadamente normal.

9.94 El Departamento de Ingeniería Civil del Virginia Tech comparó una técnica de ensayo modificada (M-5 hr) para recuperar coliformes fecales en residuos líquidos (charcos) de agua de lluvia en una área urbana con la técnica del número más probable (NMP). El departamento recolectó un total de 12 muestras de tales residuos y las analizó con las dos técnicas. Los conteos de coliformes fecales por 100 mililitros se registraron en la siguiente tabla:

Muestra	Conteo NMP	Conteo con M-5 hr
1	2300	2010
2	1200	930
3	450	400
4	210	436
5	270	4100
6	450	2090
7	154	219
8	179	169
9	192	194
10	230	174
11	340	274
12	194	183

Construya un intervalo de confianza del 90% para la diferencia entre el conteo medio de coliformes fecales que se obtuvo con la técnica M-5 hr y el que se obtuvo con la NMP. Suponga que las diferencias en los conteos se distribuyen de forma aproximadamente normal.

9.95 Se llevó a cabo un experimento para determinar si el acabado superficial tiene un efecto en el límite de resistencia a la fatiga del acero. Una teoría indica que el pulido aumenta el límite medio de resistencia a la fatiga (para la flexión inversa). Desde un punto de vista práctico, el pulido no debería tener efecto alguno sobre la desviación estándar del límite de resistencia a la fatiga, el cual se sabe, a partir de la realización de diversos

experimentos de límite de resistencia a la fatiga, que es de 4000 psi. Se realiza un experimento sobre acero al carbono al 0.4% usando especímenes sin pulido y especímenes con pulido suave. Los datos son los siguientes:

	Límite de fatiga (psi)	
	Acero al carbono al 0.4%	Acero al carbono al 0.4% sin pulir
	85,500	82,600
	91,900	82,400
	89,400	81,700
	84,000	79,500
	89,900	79,400
	78,700	69,800
	87,500	79,900
	83,100	83,400

Calcule un intervalo de confianza del 95% para la diferencia entre las medias de la población para los dos métodos. Suponga que las poblaciones se distribuyen de forma aproximadamente normal.

9.96 Un antropólogo está interesado en determinar la proporción de individuos de dos tribus indias que tienen doble remolino de cabello en la zona occipital. Suponga que toma muestras independientes de cada una de las dos tribus y encuentra que 24 de 100 individuos de la tribu A y 36 de 120 individuos de la tribu B poseen tal característica. Construya un intervalo de confianza del 95% para la diferencia $p_B - p_A$ entre las proporciones de estas dos tribus con remolinos de cabello en la zona occipital.

9.97 Un fabricante de planchas eléctricas produce estos artículos en dos plantas en las que las partes pequeñas son surtidas por el mismo proveedor. El fabricante puede ahorrar algo si le compra a un proveedor local los termostatos para la planta B. Para probar si estos nuevos termostatos son tan precisos como los anteriores le compra sólo un lote al proveedor local y los prueba en planchas a 550°F. Al final lee con un termopar las temperaturas reales y las redondea al siguiente 0.1°F más cercano. Los datos son los siguientes:

Proveedor nuevo (°F)					
530.3	559.3	549.4	544.0	551.7	566.3
549.9	556.9	536.7	558.8	538.8	543.3
559.1	555.0	538.6	551.1	565.4	554.9
550.0	554.9	554.7	536.1	569.1	
Proveedor anterior (°F)					
559.7	534.7	554.8	545.0	544.6	538.0
550.7	563.1	551.1	553.8	538.8	564.6
554.5	553.0	538.4	548.3	552.9	535.1
555.0	544.8	558.4	548.7	560.3	

Calcule un intervalo de confianza de 95% para σ_1^2/σ_2^2 y para σ_1/σ_2 , donde σ_1^2 y σ_2^2 son las varianzas de la

población de las lecturas de los termostatos del proveedor nuevo y del anterior, respectivamente.

9.98 Se afirma que la resistencia del alambre A es mayor que la del alambre B . Un experimento sobre los alambres muestra los siguientes resultados (en ohms):

Alambre A	Alambre B
0.140	0.135
0.138	0.140
0.143	0.136
0.142	0.142
0.144	0.138
0.137	0.140

Suponga varianzas iguales y explique a qué conclusiones llega si se basa en esto.

9.99 Una forma alternativa de estimación se lleva a cabo a través del método de momentos. El método consiste en igualar la media y la varianza de la población con las correspondientes media muestral \bar{x} y varianza muestral s^2 , y resolver para los parámetros; el resultado son los **estimadores por momentos**. En el caso de un solo parámetro sólo se utilizan las medias. Argumente por qué en el caso de la distribución de Poisson el estimador de máxima verosimilitud y los estimadores por momentos son iguales.

9.100 Especifique los estimadores por momentos para μ y σ^2 para la distribución normal.

9.101 Especifique los estimadores por momentos para μ y σ^2 para la distribución logarítmica normal.

9.102 Especifique los estimadores por momentos para α y β en el caso de la distribución gamma.

9.103 Se realizó una encuesta con el fin de comparar los sueldos de administradores de plantas químicas empleados en dos áreas del país: el norte y el centro-occidente. Se eligió una muestra aleatoria independiente de 300 gerentes de planta para cada una de las dos áreas. A tales gerentes se les preguntó el monto de su sueldo anual. Los resultados fueron los siguientes:

Norte	Centro-Occidente
$\bar{x}_1 = \$102,300$	$\bar{x}_2 = \$98,500$
$s_1 = \$5700$	$s_2 = \$3800$

- Construya un intervalo de confianza del 99% para $\mu_1 - \mu_2$, la diferencia en los sueldos medios.
- ¿Qué supuso en el inciso a) acerca de la distribución de los sueldos anuales para las dos áreas? ¿Es necesaria la suposición de normalidad? Explique su respuesta.
- ¿Qué supuso acerca de las dos varianzas? ¿Es razonable la suposición de igualdad de varianzas? ¿Explique!

9.104 Considere el ejercicio de repaso 9.103. Suponga que los datos aún no se han recabado. Suponga también que los estadísticos previos sugieren que $\sigma_1 = \sigma_2 = \$4000$. ¿Los tamaños de las muestras en el ejercicio de repaso 9.103 son suficientes para producir un intervalo de confianza del 95% si $\mu_1 - \mu_2$ tiene una anchura de sólo \$1000? Presente el desarrollo completo.

9.105 Un sindicato se preocupa por el notorio ausentismo de sus miembros. Los líderes del sindicato siempre habían afirmado que, en un mes típico, el 95% de sus afiliados estaban ausentes menos de 10 horas al mes. El sindicato decide verificar esto revisando una muestra aleatoria de 300 de sus miembros. Se registra el número de horas de ausencia para cada uno de los 300 miembros. Los resultados son $\bar{x} = 6.5$ horas y $s = 2.5$ horas. Utilice los datos para responder esa afirmación utilizando un límite de tolerancia unilateral y eligiendo un nivel de confianza del 99%. Asegúrese de aplicar lo que ya sabe acerca del cálculo del límite de tolerancia.

9.106 Se seleccionó una muestra aleatoria de 30 empresas que comercializan productos inalámbricos para determinar la proporción de tales empresas que implementaron software nuevo para aumentar la productividad. Resultó que 8 de las 30 empresas habían implementado tal software. Calcule un intervalo de confianza del 95% en p , la proporción verdadera de ese tipo de empresas que implementaron el nuevo software.

9.107 Remítase al ejercicio de repaso 9.106. Suponga que se desea saber si la estimación puntual $\hat{p} = 8/30$ es lo suficientemente precisa porque el intervalo de confianza alrededor de p no es tan estrecho como se requiere. Utilice \hat{p} como el estimado de p para determinar cuántas empresas habría que incluir en una muestra para obtener un intervalo de confianza del 95% con una anchura de sólo 0.05.

9.108 Un fabricante produce un artículo que se clasifica como "defectuoso" o "no defectuoso". Para estimar la proporción de productos defectuosos se tomó una muestra aleatoria de 100 artículos de la producción y se encontraron 10 defectuosos. Después de aplicar un programa de mejoramiento de la calidad se volvió a realizar el experimento. Se tomó una nueva muestra de 100 artículos y esta vez sólo 6 salieron defectuosos.

- Dado un intervalo de confianza del 95% de $p_1 - p_2$, donde p_1 y p_2 representan la proporción de artículos defectuosos de la población antes y después del mejoramiento, respectivamente.
- ¿Hay información en el intervalo de confianza que se encontró en el inciso a) que sugiera que $p_1 > p_2$? Explique su respuesta.

9.109 Se utiliza una máquina para llenar cajas de un producto en una operación de la línea de ensamble. Gran parte del interés se centra en la variabilidad del número de onzas del producto en la caja. Se sabe que la desviación estándar en el peso del producto es de 0.3 onzas. Se realizan mejoras y luego se toma una muestra aleatoria de 20 cajas, y se encuentra que la varianza de la muestra es de 0.045 onzas². Calcule un intervalo de confianza del 95% de la varianza del peso del producto. Si considera el rango del intervalo de confianza, ¿le parece que el mejoramiento en el proceso incrementó la calidad en lo que se refiere a la variabilidad? Suponga normalidad en la distribución del peso del producto.

9.110 Un grupo de consumidores está interesado en comparar los costos de operación de dos diferentes tipos de motor para automóvil. El grupo encuentra 15 propietarios cuyos automóviles tienen motor tipo *A* y 15 que tienen motor tipo *B*. Los 30 propietarios compraron sus automóviles más o menos al mismo tiempo y todos llevaron buenos registros en cierto periodo de 12 meses. Los consumidores encontraron, además, que los propietarios recorrieron aproximadamente el mismo número de millas. Los estadísticos de costo son $\bar{y}_A = \$87.00/1000$ millas, $\bar{y}_B = \$75.00/1000$ millas, $s_A = \$5.99$ y $s_B = \$4.85$. Calcule un intervalo de confianza del 95% para estimar $\mu_A - \mu_B$, la diferencia en el costo medio de operación. Suponga normalidad y varianzas iguales.

9.111 Considere el estadístico S_p^2 , el estimado agrupado de σ^2 que se estudió en la sección 9.8 y que se utiliza cuando se está dispuesto a suponer que $\sigma_1^2 = \sigma_2^2 = \sigma^2$. Demuestre que el estimador es insesgado para σ^2 [es decir, demuestre que $E(S_p^2) = \sigma^2$]. Puede utilizar los resultados de cualquier teorema o ejemplo de este capítulo.

9.112 Un grupo de investigadores del factor humano están interesados en saber cómo reaccionan los pilotos aviadores ante un estímulo dispuesto de cierta manera

en la cabina del avión. Para lograr su objetivo realizaron un experimento de simulación en un laboratorio, el cual incluyó a 15 pilotos, los que presentaron un tiempo de reacción promedio de 3.2 segundos y una desviación estándar muestral de 0.6 segundos. Resulta de interés caracterizar el extremo, es decir, el escenario del peor caso. Para conseguir esto realice lo siguiente:

- Determine un importante límite de confianza unilateral específico del 99% del tiempo medio de reacción. ¿Qué suposición, si la hubiera, debería hacer acerca de la distribución de los tiempos de reacción?
- Determine un intervalo unilateral de predicción del 99% e interprete su significado. ¿Debería usted suponer algo sobre la distribución de los tiempos de reacción para calcular este límite?
- Calcule un límite de tolerancia unilateral con una confianza del 99% que incluya al 95% de los tiempos de reacción. Nuevamente, de ser necesario, interprete o suponga algo acerca de la distribución. [Nota: Los valores del límite de tolerancia unilateral también se incluyen en la tabla A.7].

9.113 Cierta proveedor fabrica un tipo de tapete de hule que vende a las empresas automotrices. El material que utiliza para los tapetes debe tener ciertas características de dureza. Ocasionalmente detecta tapetes defectuosos en el proceso y los rechaza. El proveedor afirma que la proporción de tapetes defectuosos es de 0.05, pero como un cliente que compró los tapetes desafió su afirmación, realizó un experimento en el que se probaron 400 tapetes y se encontraron 17 defectuosos.

- Calcule un intervalo de confianza bilateral del 95% de la proporción de tapetes defectuosos.
- Calcule un intervalo de confianza unilateral del 95% adecuado de la proporción de tapetes defectuosos.
- Interprete los intervalos de ambos incisos y comente acerca de la afirmación hecha por el proveedor.

9.15 Posibles riesgos y errores conceptuales: relación con el material de otros capítulos

El concepto de *intervalo de confianza de muestra grande* en una población a menudo confunde a los alumnos principiantes. Se basa en la idea de que incluso cuando se desconoce σ y no se está convencido de que la distribución que se muestrea es normal, se puede calcular un intervalo de confianza para μ a partir de

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$$

En la práctica es común que se utilice esta fórmula cuando la muestra es demasiado pequeña. El origen de este intervalo de muestra grande es, por supuesto, el teorema del

límite central (TLC), con el cual la normalidad no es necesaria. Aquí el TLC requiere una σ conocida, de la cual s sólo es un estimado. Por lo tanto, n debe ser al menos tan grande como 30 y la distribución subyacente debe tener una simetría similar, en cuyo caso el intervalo sigue siendo una aproximación.

Hay casos en que la aplicación práctica del material de este capítulo depende en gran medida del contexto específico. Un ejemplo muy importante es el uso de la distribución t para el intervalo de confianza de μ cuando se desconoce σ . En términos estrictos, el uso de la distribución t requiere que la distribución de donde se toma la muestra sea normal. Sin embargo, es bien sabido que cualquier aplicación de la distribución t es razonablemente insensible, es decir, **robusta**, a la suposición de normalidad. Esto representa una de esas situaciones afortunadas que ocurren con frecuencia en el campo de la estadística, donde no se sostiene un supuesto básico y "¡todo resulta bien!" Sin embargo, la población de la que se toma la muestra no se puede desviar mucho de la normalidad. Por consiguiente, a menudo se recurrirá a las gráficas de probabilidad normal estudiadas en el capítulo 8 y las pruebas de bondad del ajuste que se presentarán en el capítulo 10 para atribuir algún sentido de "cercanía a la normalidad". Esta idea de "robustez a la normalidad" se volverá a presentar en el capítulo 10.

Por experiencia sabemos que uno de los más graves "usos incorrectos de la estadística" en la práctica surge de la confusión sobre las diferencias en la interpretación de los tipos de intervalos estadísticos. Por consiguiente, la subsección de este capítulo en la que se examinan las diferencias entre los tres tipos de intervalos es importante. Es muy probable que en la práctica se utilice en exceso el intervalo de confianza, es decir, que se emplee cuando no es la media lo que interesa en realidad, sino la cuestión de: "¿en dónde va a caer la siguiente observación?", o la a menudo más importante cuestión de: "¿en dónde se ubica la mayor parte de la distribución?" Éstas son preguntas fundamentales que no se pueden responder calculando un intervalo de la media. A menudo resulta confusa la interpretación de un intervalo de confianza. Es tentador concluir que hay una probabilidad de 0.95 de que el parámetro caiga dentro del intervalo. Aunque se trata de una interpretación correcta del **intervalo posterior bayesiano** (para mayores referencias sobre la inferencia bayesiana véase el capítulo 18), no es una interpretación adecuada de la frecuencia.

El intervalo de confianza tan sólo sugiere que si se realiza el experimento y los datos se observan una y otra vez, aproximadamente 95% de tales intervalos contendrá el parámetro verdadero. Cualquier alumno principiante de la estadística práctica debería tener muy claras las diferencias entre estos intervalos estadísticos.

Otro posible y grave uso incorrecto de la estadística es el que se cometería si se aplicara la distribución χ^2 a un intervalo de confianza de una sola varianza. De nuevo, se supone normalidad en la distribución de donde se toma la muestra. A diferencia del resultado de utilizar la distribución t , la prueba χ^2 para esta aplicación **no es robusta para la suposición de normalidad** (esto significa que cuando la distribución subyacente no es normal, la distribución muestral de $\frac{(n-1)S^2}{\sigma^2}$ se aparta mucho de χ^2). En consecuencia, el uso estricto de la prueba de bondad de ajuste (véase el capítulo 10) y de las gráficas de probabilidad normal, o de la prueba y las gráficas, puede ser muy importante en esos contextos. En los siguientes capítulos se proporcionará más información sobre este tema general.

Capítulo 10

Pruebas de hipótesis de una y dos muestras

10.1 Hipótesis estadísticas: conceptos generales

Como se expuso en el capítulo 9, a menudo el problema al que se enfrentan el científico o el ingeniero no es tanto la estimación de un parámetro de la población, sino la formación de un procedimiento de decisión que se base en los datos y que pueda producir una conclusión acerca de algún sistema científico. Por ejemplo, un investigador médico puede decidir con base en evidencia experimental si beber café incrementa el riesgo de cáncer en los seres humanos; un ingeniero quizá tenga que decidir con base en datos muestrales si hay una diferencia entre la precisión de un tipo de medidor y la de otro; o tal vez un sociólogo desee reunir los datos apropiados que le permitan decidir si el tipo de sangre y el color de ojos de un individuo son variables independientes. En cada uno de estos casos el científico o el ingeniero *postulan o conjeturan* algo acerca de un sistema. Además, cada uno debe utilizar datos experimentales y tomar decisiones basadas en ellos. En cada caso la conjetura se puede expresar en forma de hipótesis estadística. Los procedimientos que conducen a la aceptación o al rechazo de hipótesis estadísticas como éstas comprenden una área importante de la inferencia estadística. Empecemos por definir con precisión lo que entendemos por **hipótesis estadística**.

Definición 10.1: Una **hipótesis estadística** es una aseveración o conjetura respecto a una o más poblaciones.

La verdad o falsedad de una hipótesis estadística nunca se sabe con absoluta certeza, a menos que se examine toda la población, lo cual, por supuesto, sería poco práctico en la mayoría de las situaciones. En vez de eso se toma una muestra aleatoria de la población de interés y se utilizan los datos contenidos en ella para proporcionar evidencia que respalde o no la hipótesis. La evidencia de la muestra que es inconsistente con la hipótesis planteada conduce al rechazo de la misma.

El papel que desempeña la probabilidad en la prueba de hipótesis

Debería quedar claro al lector que un procedimiento de toma de decisiones debe implicar la conciencia de la *probabilidad de llegar a una conclusión errónea*. Por ejemplo, suponga que la hipótesis que postuló el ingeniero es que la fracción p de artículos defectuosos en cierto proceso es 0.10. El experimento consiste en observar una muestra aleatoria del producto en cuestión. Suponga que se prueban 100 artículos y que se encuentran 12 defectuosos. Es razonable concluir que esta evidencia no rechaza la condición de que el parámetro binomial $p = 0.10$, por lo que puede provocar que no se rechace la hipótesis. Sin embargo, también puede provocar que no se refute $p = 0.12$, o quizá incluso $p = 0.15$. Como resultado, el lector se debe acostumbrar a la idea de que **el rechazo de una hipótesis implica que fue refutada por la evidencia de la muestra**. En otras palabras, **el rechazo significa que existe una pequeña probabilidad de obtener la información muestral observada cuando, de hecho, la hipótesis es verdadera**. Por ejemplo, en la hipótesis de la proporción de artículos defectuosos, una muestra de 100 artículos que revela que hay 20 defectuosos es ciertamente evidencia para el rechazo. ¿Por qué? Si en realidad $p = 0.10$, la probabilidad de obtener 20 o más artículos defectuosos es aproximadamente de 0.002. Con el pequeño riesgo resultante de llegar a una conclusión errónea parecería seguro **rechazar la hipótesis** de que $p = 0.10$. En otras palabras, el rechazo de una hipótesis tiende a casi “descartar” la hipótesis. Por otro lado, es muy importante enfatizar que la aceptación o, más bien, la falta de rechazo no descarta otras posibilidades. Como resultado, *el analista de datos establece una conclusión firme cuando se rechaza una hipótesis*.

En el planteamiento formal de una hipótesis a menudo influye la estructura de la probabilidad de una conclusión errónea. Si el científico está interesado en *apoyar firmemente* un argumento, espera llegar a éste en la forma del rechazo de una hipótesis. Si el investigador médico desea mostrar evidencia sólida a favor del argumento de que beber café aumenta el riesgo de contraer cáncer, la hipótesis a probar debería tener la forma “el riesgo de desarrollar cáncer no aumenta como consecuencia de beber café”. Como resultado, el argumento se obtiene mediante un rechazo. De manera similar, para apoyar la afirmación de que un tipo de medidores es más preciso que otro, el ingeniero prueba la hipótesis de que no hay diferencia en la precisión de los dos tipos de medidores.

Lo anterior implica que cuando el analista de datos formaliza la evidencia experimental con base en la prueba de hipótesis, es muy importante el **planteamiento formal de la hipótesis**.

La hipótesis nula y la hipótesis alternativa

La estructura de la prueba de hipótesis se establece usando el término **hipótesis nula**, el cual se refiere a cualquier hipótesis que se desea probar y se denota con H_0 . El rechazo de H_0 conduce a la aceptación de una **hipótesis alternativa**, que se denota con H_1 . La comprensión de las diferentes funciones que desempeñan la hipótesis nula (H_0) y la hipótesis alternativa (H_1) es fundamental para entender los principios de la prueba de hipótesis. La hipótesis alternativa H_1 por lo general representa la *pregunta que se responderá o la teoría que se probará*, por lo que su especificación es muy importante. La hipótesis nula H_0 *anula o se opone a H_1* y a menudo es el complemento lógico de H_1 . A medida que el lector aprenda más sobre la prueba de hipótesis notará que el analista llega a una de las siguientes dos conclusiones:

rechazar H_0 a favor de H_1 debido a evidencia suficiente en los datos o
no rechazar H_0 debido a evidencia insuficiente en los datos.

Observe que las conclusiones no implican una "aceptación de H_0 " formal y literal. La aseveración de H_0 a menudo representa el "status quo" contrario a una nueva idea, conjetura, etcétera, enunciada en H_1 ; en tanto que no rechazar H_0 representa la conclusión adecuada. En nuestro ejemplo binomial la cuestión práctica podría ser el interés en que la probabilidad histórica de artículos defectuosos de 0.10 ya no sea verdadera. De hecho, la conjetura podría ser que p excede a 0.10. Entonces podríamos afirmar que

$$H_0: p = 0.10,$$

$$H_1: p > 0.10.$$

Ahora, 12 artículos defectuosos de cada 100 no refutan $p = 0.10$, por lo que la conclusión es "no rechazar H_0 ". Sin embargo, si los datos revelan 20 artículos defectuosos de cada 100, la conclusión sería "rechazar H_0 " a favor de H_1 ; $p > 0.10$.

Aunque las aplicaciones de la prueba de hipótesis son muy abundantes en trabajos científicos y de ingeniería, quizás el mejor ejemplo para un principiante sea el dilema que enfrenta el jurado en un juicio. Las hipótesis nula y alternativa son

H_0 : el acusado es inocente,

H_1 : el acusado es culpable.

La acusación proviene de una sospecha de culpabilidad. La hipótesis H_0 (el status quo) se establece en oposición a H_1 y se mantiene a menos que se respalde H_1 con evidencia "más allá de una duda razonable". Sin embargo, en este caso "no rechazar H_0 " no implica inocencia, sino sólo que la evidencia fue insuficiente para lograr una condena. Por lo tanto, el jurado no necesariamente *acepta H_0* sino que *no rechaza H_0* .

10.2 Prueba de una hipótesis estadística

Para ilustrar los conceptos que se utilizan al probar una hipótesis estadística acerca de una población considere el siguiente ejemplo. Se sabe que, después de un periodo de dos años, cierto tipo de vacuna contra un virus que produce resfriado ya sólo es 25% eficaz. Suponga que se eligen 20 personas al azar y se les aplica una vacuna nueva, un poco más costosa, para determinar si protege contra el mismo virus durante un periodo más largo. (En un estudio real de este tipo el número de participantes que reciben la nueva vacuna podría ascender a varios miles. Aquí la muestra es de 20 sólo porque lo único que se busca es demostrar los pasos básicos para realizar una prueba estadística). Si más de 8 individuos de los que reciben la nueva vacuna superan el lapso de 2 años sin contraer el virus, la nueva vacuna se considerará superior a la que se usa en la actualidad. El requisito de que el número exceda a 8 es algo arbitrario, aunque parece razonable, ya que representa una mejoría modesta sobre las 5 personas que se esperaría recibieran protección si fueran inoculadas con la vacuna que actualmente está en uso. En esencia probamos la hipótesis nula de que la nueva vacuna es igual de eficaz después de un periodo de 2 años que la que se utiliza en la actualidad. La hipótesis alternativa es que la nueva vacuna es

mejor, y esto equivale a poner a prueba la hipótesis de que el parámetro binomial para la probabilidad de un éxito en un ensayo dado es $p = 1/4$, contra la alternativa de que $p > 1/4$. Esto por lo general se escribe como se indica a continuación:

$$\begin{aligned} H_0: p &= 0.25, \\ H_1: p &> 0.25. \end{aligned}$$

El estadístico de prueba

El estadístico de prueba en el cual se basa nuestra decisión es X , el número de individuos en nuestro grupo de prueba que reciben protección de la nueva vacuna durante un periodo de al menos 2 años. Los valores posibles de X , de 0 a 20, se dividen en dos grupos: los números menores o iguales que 8 y aquellos mayores que 8. Todos los posibles valores mayores que 8 constituyen la **región crítica**. El último número que observamos al pasar a la región crítica se llama **valor crítico**. En nuestro ejemplo el valor crítico es el número 8. Por lo tanto, si $x > 8$, rechazamos H_0 a favor de la hipótesis alternativa H_1 . Si $x \leq 8$, no rechazamos H_0 . Este criterio de decisión se ilustra en la figura 10.1.

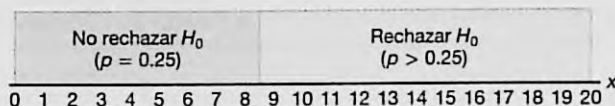


Figura 10.1: Criterio de decisión para probar $p = 0.25$ contra $p > 0.25$.

La probabilidad de un error tipo 1

El procedimiento de toma de decisiones recién descrito podría conducir a cualquiera de dos conclusiones erróneas. Por ejemplo, es probable que la nueva vacuna no sea mejor que la que se usa en la actualidad (H_0 verdadera) y, sin embargo, en este grupo específico de individuos seleccionados aleatoriamente más de 8 pasan el periodo de 2 años sin contraer el virus. Si rechazáramos H_0 a favor de H_1 cuando, de hecho, H_0 es verdadera, cometeríamos un error que se conoce como **error tipo I**.

Definición 10.2: El rechazo de la hipótesis nula cuando es verdadera se denomina **error tipo I**.

Si 8 o menos miembros del grupo superan exitosamente el periodo de 2 años y no concluimos que la nueva vacuna es mejor cuando en realidad sí lo es (H_1 verdadera), cometemos un segundo tipo de error, el de no rechazar la hipótesis H_0 cuando en realidad es falsa. A este error se le conoce como **error tipo II**.

Definición 10.3: No rechazar la hipótesis nula cuando es falsa se denomina **error tipo II**.

Al probar cualquier hipótesis estadística, hay cuatro situaciones posibles que determinan si nuestra decisión es correcta o errónea. Estas cuatro situaciones se resumen en

la tabla 10.1.

Tabla 10.1: Situaciones posibles al probar una hipótesis estadística.

	H_0 es verdadera	H_0 es falsa
No rechazar H_0	Decisión correcta	Error tipo II
Rechazar H_0	Error tipo I	Decisión correcta

La probabilidad de cometer un error tipo I, también llamada **nivel de significancia**, se denota con la letra griega α . En nuestro ejemplo un error tipo I ocurriría si más de 8 individuos inoculados con la nueva vacuna superan el periodo de 2 años sin contraer el virus y los investigadores concluyen que la nueva vacuna es mejor, cuando en realidad es igual a la vacuna que se utiliza en la actualidad. Por lo tanto, si X es el número de individuos que permanecen sin contraer el virus por al menos dos años,

$$\begin{aligned}\alpha = P(\text{error tipo I}) &= P\left(X > 8 \text{ cuando } p = \frac{1}{4}\right) = \sum_{x=9}^{20} b\left(x; 20, \frac{1}{4}\right) \\ &= 1 - \sum_{x=0}^8 b\left(x; 20, \frac{1}{4}\right) = 1 - 0.9591 = 0.0409.\end{aligned}$$

Decimos que la hipótesis nula, $p = 1/4$, se prueba al nivel de significancia $\alpha = 0.0409$. En ocasiones el nivel de significancia se conoce como **tamaño de la prueba**. Una región crítica de tamaño 0.0409 es muy pequeña y, por lo tanto, es poco probable que se cometa un error de tipo I. En consecuencia, sería poco probable que más de 8 individuos permanecieran inmunes a un virus durante 2 años utilizando una vacuna nueva que en esencia es equivalente a la que actualmente está en el mercado.

La probabilidad de un error tipo II

La probabilidad de cometer un error tipo II, que se denota con β , es imposible de calcular a menos que tengamos una hipótesis alternativa específica. Si probamos la hipótesis nula $p = 1/4$ contra la hipótesis alternativa $p = 1/2$, entonces podremos calcular la probabilidad de no rechazar H_0 cuando es falsa. Simplemente calculamos la probabilidad de obtener 8 o menos en el grupo que supera el periodo de 2 años cuando $p = 1/2$. En este caso,

$$\begin{aligned}\beta = P(\text{error tipo II}) &= P\left(X \leq 8 \text{ cuando } p = \frac{1}{2}\right) \\ &= \sum_{x=0}^8 b\left(x; 20, \frac{1}{2}\right) = 0.2517.\end{aligned}$$

Se trata de una probabilidad elevada que indica un procedimiento de prueba en el cual es muy probable que se rechace la nueva vacuna cuando, de hecho, es mejor a la que está actualmente en uso. De manera ideal, es preferible utilizar un procedimiento de prueba con el cual haya pocas probabilidades de cometer el error tipo I y el error tipo II.

Es posible que el director del programa de prueba esté dispuesto a cometer un error tipo II si la vacuna más costosa no es significativamente mejor. De hecho, la única

ocasión en la que desea evitar un error tipo II es cuando el verdadero valor de p es de al menos 0.7. Si $p = 0.7$, este procedimiento de prueba da

$$\begin{aligned}\beta &= P(\text{error tipo II}) = P(X \leq 8 \text{ cuando } p = 0.7) \\ &= \sum_{x=0}^8 b(x; 20, 0.7) = 0.0051.\end{aligned}$$

Con una probabilidad tan pequeña de cometer un error tipo II es muy improbable que se rechace la nueva vacuna cuando tiene una efectividad de 70% después de un periodo de 2 años. A medida que la hipótesis alternativa se aproxima a la unidad, el valor de β tiende a disminuir hasta cero.

El papel que desempeñan α , β y el tamaño de la muestra

Supongamos que el director del programa de prueba no está dispuesto a cometer un error tipo II cuando la hipótesis alternativa $p = 1/2$ es verdadera, aun cuando se encuentre que la probabilidad de tal error es $\beta = 0.2517$. Siempre es posible reducir β aumentando el tamaño de la región crítica. Por ejemplo, considere lo que les sucede a los valores de α y β cuando cambiamos nuestro valor crítico a 7, de manera que todos los valores mayores que 7 caigan en la región crítica y aquellos menores o iguales que 7 caigan en la región de no rechazo. Así, al probar $p = 1/4$ contra la hipótesis alternativa $p = 1/2$, encontramos que

$$\begin{aligned}\alpha &= \sum_{x=8}^{20} b\left(x; 20, \frac{1}{4}\right) = 1 - \sum_{x=0}^7 b\left(x; 20, \frac{1}{4}\right) = 1 - 0.8982 = 0.1018 \\ \beta &= \sum_{x=0}^7 b\left(x; 20, \frac{1}{2}\right) = 0.1316.\end{aligned}$$

Al adoptar un nuevo procedimiento de toma de decisiones, reducimos la probabilidad de cometer un error tipo II a costa de aumentar la probabilidad de cometer un error tipo I. Para un tamaño muestral fijo, una disminución en la probabilidad de un error por lo general tendrá como resultado un incremento en la probabilidad del otro error. Por fortuna, **la probabilidad de cometer ambos tipos de errores se puede reducir aumentando el tamaño de la muestra.** Considere el mismo problema usando una muestra aleatoria de 100 individuos. Si más de 36 miembros del grupo superan el periodo de 2 años, rechazamos la hipótesis nula de $p = 1/4$ y aceptamos la hipótesis alternativa de $p > 1/4$. El valor crítico ahora es 36. Todos los valores posibles mayores de 36 constituyen la región crítica y todos los valores posibles menores o iguales que 36 caen en la región de aceptación.

Para determinar la probabilidad de cometer un error tipo I debemos utilizar la aproximación a la curva normal con

$$\mu = np = (100) \left(\frac{1}{4}\right) = 25 \quad \text{y} \quad \sigma = \sqrt{npq} = \sqrt{(100)(1/4)(3/4)} = 4.33.$$

Con respecto a la figura 10.2, necesitamos el área bajo la curva normal a la derecha de $x = 36.5$. El valor z correspondiente es

$$z = \frac{36.5 - 25}{4.33} = 2.66.$$

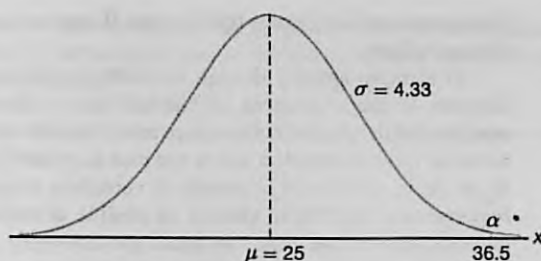


Figura 10.2: Probabilidad de un error tipo I.

En la tabla A.3 encontramos que

$$\begin{aligned}\alpha &= P(\text{error tipo I}) = P\left(X > 36 \text{ cuando } p = \frac{1}{4}\right) \approx P(Z > 2.66) \\ &= 1 - P(Z < 2.66) = 1 - 0.9961 = 0.0039.\end{aligned}$$

Si H_0 es falsa y el verdadero valor de H_1 es $p = 1/2$, determinamos la probabilidad de un error tipo II usando la aproximación a la curva normal con

$$\mu = np = (100)(1/2) = 50 \quad \text{y} \quad \sigma = \sqrt{npq} = \sqrt{(100)(1/2)(1/2)} = 5.$$

La probabilidad de que un valor caiga en la región de no rechazo cuando H_0 es verdadera es dada por el área de la región sombreada a la izquierda de $x = 36.5$ en la figura 10.3. El valor z que corresponde a $x = 36.5$ es

$$z = \frac{36.5 - 50}{5} = -2.7.$$

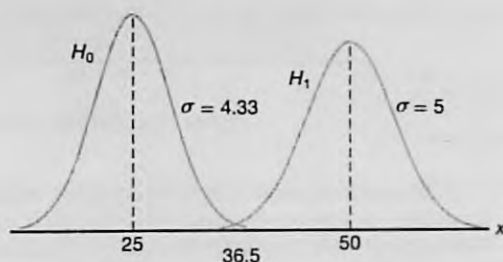


Figura 10.3: Probabilidad de un error tipo II.

Por lo tanto,

$$\beta = P(\text{error tipo II}) = P\left(X \leq 36 \text{ cuando } p = \frac{1}{2}\right) \approx P(Z < -2.7) = 0.0035.$$

Evidentemente, los errores tipo I y tipo II rara vez ocurren si el experimento consta de 100 individuos.

El ejemplo anterior destaca la estrategia del científico en la prueba de hipótesis. Después de que se plantean las hipótesis nula y alternativa es importante considerar la sensibilidad del procedimiento de prueba. Con esto queremos decir que debería determinarse un valor razonable a una α fija para la probabilidad de aceptar de manera errónea H_0 , es decir, el valor de β , cuando la verdadera situación representa alguna *desviación importante de H_0* . Por lo general, es posible determinar un valor para el tamaño de la muestra, para el que existe un equilibrio razonable entre los valores de α y β que se calcula de esta manera. El problema de la vacuna es un ejemplo.

Ilustración con una variable aleatoria continua

Los conceptos que se analizan aquí para una población discreta también se pueden aplicar a variables aleatorias continuas. Considere la hipótesis nula de que el peso promedio de estudiantes hombres en cierta universidad es de 68 kilogramos, contra la hipótesis alternativa de que es diferente a 68. Es decir, deseamos probar

$$\begin{aligned} H_0: \mu &= 68, \\ H_1: \mu &\neq 68. \end{aligned}$$

La hipótesis alternativa nos permite la posibilidad de que $\mu < 68$ o $\mu > 68$.

Una media muestral que caiga cerca del valor hipotético de 68 se consideraría como evidencia a favor de H_0 . Por otro lado, una media muestral considerablemente menor que o mayor que 68 sería evidencia en contra de H_0 y, por lo tanto, favorecería a H_1 . La media muestral es el estadístico de prueba en este caso. Una región crítica para el estadístico de prueba se puede elegir de manera arbitraria como los dos intervalos $\bar{x} < 67$ y $\bar{x} > 69$. La región de no rechazo será entonces el intervalo $67 \leq \bar{x} \leq 69$. Este criterio de decisión se ilustra en la figura 10.4.

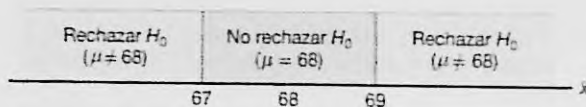


Figura 10.4: Región crítica (en azul).

Utilicemos ahora el criterio de decisión de la figura 10.4 para calcular las probabilidades de cometer los errores tipo I y tipo II cuando probemos la hipótesis nula $\mu = 68$ kilogramos contra la alternativa $\mu \neq 68$ kilogramos.

Suponga que la desviación estándar de la población de pesos es $\sigma = 3.6$. Para muestras grandes podemos sustituir s por σ si no disponemos de ninguna otra estimación de σ . Nuestro estadístico de decisión, que se basa en una muestra aleatoria de tamaño $n = 36$, será \bar{X} , el estimador más eficaz de μ . Del teorema del límite central sabemos que la distribución muestral de \bar{X} es aproximadamente normal con desviación estándar $\sigma_{\bar{X}} = \sigma/\sqrt{n} = 3.6/6 = 0.6$.

La probabilidad de cometer un error tipo I, o el nivel de significancia de nuestra prueba, es igual a la suma de las áreas sombreadas en cada cola de la distribución en la figura 10.5. Por lo tanto,

$$\alpha = P(\bar{X} < 67 \text{ cuando } \mu = 68) + P(\bar{X} > 69 \text{ cuando } \mu = 68).$$

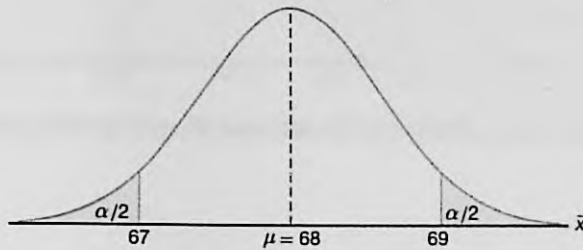


Figura 10.5: Región crítica para probar $\mu = 68$ contra $\mu \neq 68$.

Los valores z correspondientes a $\bar{x}_1 = 67$ y $\bar{x}_2 = 69$ cuando H_0 es verdadera son

$$z_1 = \frac{67 - 68}{0.6} = -1.67 \quad \text{y} \quad z_2 = \frac{69 - 68}{0.6} = 1.67.$$

Por lo tanto,

$$\alpha = P(Z < -1.67) + P(Z > 1.67) = 2P(Z < -1.67) = 0.0950.$$

Por consiguiente, 9.5% de todas las muestras de tamaño 36 nos conducirían a rechazar $\mu = 68$ kilogramos cuando, de hecho, ésta es verdadera. Para reducir α tenemos que elegir entre aumentar el tamaño de la muestra o ampliar la región de no rechazo. Suponga que aumentamos el tamaño de la muestra a $n = 64$. Entonces $\sigma_{\bar{x}} = 3.6/8 = 0.45$. En consecuencia,

$$z_1 = \frac{67 - 68}{0.45} = -2.22 \quad \text{y} \quad z_2 = \frac{69 - 68}{0.45} = 2.22.$$

Por lo tanto,

$$\alpha = P(Z < -2.22) + P(Z > 2.22) = 2P(Z < -2.22) = 0.0264.$$

La reducción de α no es suficiente por sí misma para garantizar un buen procedimiento de prueba. Debemos evaluar β para varias hipótesis alternativas. Si es importante rechazar H_0 cuando la media verdadera sea algún valor $\mu \geq 70$ o $\mu \leq 66$, entonces se debería calcular y examinar la probabilidad de cometer un error tipo II para las alternativas $\mu = 66$ y $\mu = 70$. Debido a la simetría, sólo es necesario considerar la probabilidad de no rechazar la hipótesis nula $\mu = 68$ cuando la alternativa $\mu = 70$ es verdadera. Cuando la media muestral \bar{x} caiga entre 67 y 69, cuando H_1 sea verdadera, resultará un error tipo II. Por lo tanto, remitiéndonos a la figura 10.6 encontramos que

$$\beta = P(67 \leq \bar{X} \leq 69 \text{ cuando } \mu = 70).$$

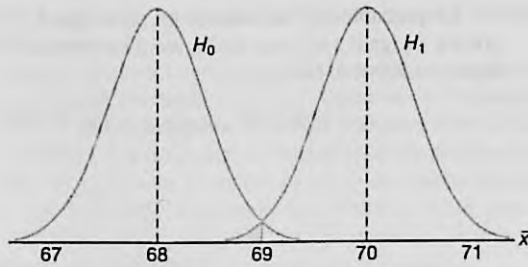


Figura 10.6: Probabilidad del error tipo II al probar $\mu = 68$ contra $\mu = 70$.

Los valores z que corresponden a $\bar{x}_1 = 67$ y $\bar{x}_2 = 69$ cuando H_1 es verdadera son

$$z_1 = \frac{67 - 70}{0.45} = -6.67 \quad \text{y} \quad z_2 = \frac{69 - 70}{0.45} = -2.22.$$

Por lo tanto,

$$\begin{aligned} \beta &= P(-6.67 < Z < -2.22) = P(Z < -2.22) - P(Z < -6.67) \\ &= 0.0132 - 0.0000 = 0.0132. \end{aligned}$$

Si el valor verdadero de μ es la alternativa $\mu = 66$, el valor de β nuevamente será 0.0132. Para todos los valores posibles de $\mu < 66$ o $\mu > 70$, el valor de β será incluso más pequeño cuando $n = 64$ y, en consecuencia, habrá poca oportunidad de no rechazar H_0 cuando sea falsa.

La probabilidad de cometer un error tipo II aumenta rápidamente cuando el valor verdadero de μ se aproxima al valor hipotético pero no es igual a éste. Desde luego, ésta suele ser la situación en la que no nos importa cometer un error tipo II. Por ejemplo, si la hipótesis alternativa $\mu = 68.5$ es verdadera, no nos importa cometer un error tipo II al concluir que la respuesta verdadera es $\mu = 68$. La probabilidad de cometer tal error será elevada cuando $n = 64$. Al remitirnos a la figura 10.7, tenemos

$$\beta = P(67 \leq \bar{X} \leq 69 \text{ cuando } \mu = 68.5).$$

Los valores z correspondientes a $\bar{x}_1 = 67$ y $\bar{x}_2 = 69$ cuando $\mu = 68.5$ son

$$z_1 = \frac{67 - 68.5}{0.45} = -3.33 \quad \text{y} \quad z_2 = \frac{69 - 68.5}{0.45} = 1.11.$$

Por lo tanto,

$$\begin{aligned} \beta &= P(-3.33 < Z < 1.11) = P(Z < 1.11) - P(Z < -3.33) \\ &= 0.8665 - 0.0004 = 0.8661. \end{aligned}$$

Los ejemplos anteriores ilustran las siguientes propiedades importantes:

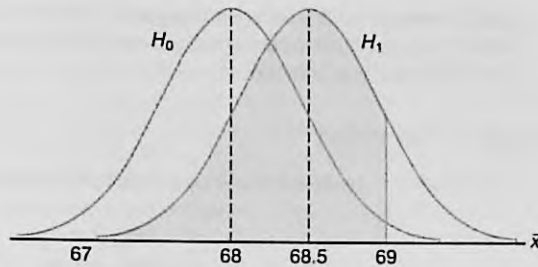


Figura 10.7: Error tipo II para la prueba de $\mu = 68$ contra $\mu = 68.5$.

Propiedades importantes de una prueba de hipótesis

1. Los errores tipo I y tipo II están relacionados. Por lo general una disminución en la probabilidad de cometer uno da como resultado un incremento en la probabilidad de cometer el otro.
2. El tamaño de la región crítica y, por lo tanto, la probabilidad de cometer un error tipo I, siempre se puede reducir ajustando el (los) valor(es) crítico(s).
3. Un aumento en el tamaño de la muestra n reducirá α y β de forma simultánea.
4. Si la hipótesis nula es falsa, β es un máximo cuando el valor verdadero de un parámetro se aproxima al valor hipotético. Cuanto más grande sea la distancia entre el valor verdadero y el valor hipotético, más pequeña será β .

Definición 10.4: La **potencia** de una prueba es la probabilidad de rechazar H_0 dado que una alternativa específica es verdadera.

La potencia de una prueba se puede calcular como $1 - \beta$. A menudo **diferentes tipos de pruebas se comparan contrastando propiedades de potencia**. Considere el caso anterior en el que probamos $H_0: \mu = 68$ y $H_1: \mu \neq 68$. Como antes, suponga que nos interesa evaluar la sensibilidad de la prueba, la cual es determinada por la regla de que no rechazamos H_0 si $67 \leq \bar{x} \leq 69$. Buscamos la capacidad de la prueba para rechazar H_0 de manera adecuada cuando en realidad $\mu = 68.5$. Vimos que la probabilidad de un error tipo II es dada por $\beta = 0.8661$. Por consiguiente, la **potencia** de la prueba es $1 - 0.8661 = 0.1339$. En cierto sentido, la potencia es una medida más sucinta de cuán sensible es la prueba para detectar diferencias entre una media de 68 y otra de 68.5. En este caso, si μ es verdaderamente 68.5, la prueba como se describe *rechazará de forma adecuada H_0 sólo 13.39% de las veces*. Como resultado, la prueba no sería buena si es importante que el analista tenga una oportunidad razonable de distinguir realmente entre una media de 68.0 (que especifica H_0) y una media de 68.5. De lo anterior resulta claro que para producir una potencia deseable, digamos, mayor que 0.8, es necesario incrementar α o aumentar el tamaño de la muestra.

Hasta ahora gran parte del análisis de la prueba de hipótesis se ha enfocado en los principios y las definiciones. En las secciones que siguen seremos más específicos y

clasificaremos las hipótesis en categorías. También estudiaremos pruebas de hipótesis sobre varios parámetros de interés. Comenzamos estableciendo la diferencia entre hipótesis unilaterales y bilaterales.

Pruebas de una y dos colas

Una prueba de cualquier hipótesis estadística donde la alternativa es **unilateral**, como

$$H_0: \theta = \theta_0,$$

$$H_1: \theta > \theta_0,$$

o quizás

$$H_0: \theta = \theta_0,$$

$$H_1: \theta < \theta_0,$$

se denomina **prueba de una sola cola**. Anteriormente en esta sección se hizo referencia al **estadístico de prueba** para una hipótesis. Por lo general la región crítica para la hipótesis alternativa $\theta > \theta_0$ yace en la cola derecha de la distribución del estadístico de prueba, en tanto que la región crítica para la hipótesis alternativa $\theta < \theta_0$ yace por completo en la cola izquierda. (En cierto sentido el símbolo de desigualdad señala la dirección en donde se encuentra la región crítica). En el experimento de la vacuna se utilizó una prueba de una sola cola para probar la hipótesis $p = 1/4$ contra la alternativa unilateral $p > 1/4$ para la distribución binomial. La región crítica de una sola cola por lo general es evidente; el lector debería visualizar el comportamiento del estadístico de prueba y observar la *señal* evidente que produciría evidencia que respalde la hipótesis alternativa.

La prueba de cualquier hipótesis alternativa donde la alternativa es **bilateral**, como

$$H_0: \theta = \theta_0,$$

$$H_1: \theta \neq \theta_0,$$

se denomina **prueba de dos colas**, ya que la región crítica se divide en dos partes, a menudo con probabilidades iguales en cada cola de la distribución del estadístico de prueba. La hipótesis alternativa $\theta \neq \theta_0$ establece que $\theta < \theta_0$ o que $\theta > \theta_0$. Se utilizó una prueba de dos colas para probar la hipótesis nula $\mu = 68$ kilogramos contra la alternativa bilateral $\mu \neq 68$ kilogramos en el ejemplo de la población continua de los pesos de estudiantes.

¿Cómo se eligen las hipótesis nula y alternativa?

Con frecuencia la hipótesis nula H_0 se plantea usando el *signo de igualdad*. Con este método se observa claramente cómo se controla la probabilidad de cometer un error tipo I. Sin embargo, hay situaciones en que “no rechazar H_0 ” implica que el parámetro θ podría ser cualquier valor definido por el complemento natural de la hipótesis alternativa. Por ejemplo, en el caso de la vacuna, donde la hipótesis alternativa es $H_1: p > 1/4$, es muy posible que el no rechazo de H_0 no pueda descartar un valor de p menor que $1/4$. Sin embargo, es evidente que en el caso de las pruebas de una cola la consideración más importante es el planteamiento de la alternativa.

La decisión de plantear una prueba de una cola o una de dos colas depende de la conclusión que se obtenga si se rechaza H_0 . La ubicación de la región crítica sólo se puede determinar después de que se plantea H_1 . Por ejemplo, al probar una medicina nueva se establece la hipótesis de que no es mejor que las medicinas similares que actualmente hay en el mercado y se prueba contra la hipótesis alternativa de que la medicina nueva es mejor. Esta hipótesis alternativa dará como resultado una prueba de una sola cola, con la región crítica en la cola derecha. Sin embargo, si deseamos comparar una nueva técnica de enseñanza con el procedimiento convencional del salón de clases, la hipótesis alternativa debe permitir que el nuevo método sea inferior o superior al procedimiento convencional. Por lo tanto, la prueba sería de dos colas con la región crítica dividida en partes iguales, de manera que caiga en los extremos de las colas izquierda y derecha de la distribución de nuestro estadístico.

Ejemplo 10.1: Un fabricante de cierta marca de cereal de arroz afirma que el contenido promedio de grasa saturada no excede a 1.5 gramos por porción. Plantee las hipótesis nula y alternativa que se utilizarán para probar esta afirmación y establezca en dónde se localiza la región crítica.

Solución: La afirmación del fabricante se rechazará sólo si μ es mayor que 1.5 miligramos y no se rechazará si μ es menor o igual que 1.5 miligramos. Entonces, probamos

$$H_0: \mu = 1.5,$$

$$H_1: \mu > 1.5.$$

El hecho de no rechazar H_0 no descarta valores menores que 1.5 miligramos. Como tenemos una prueba de una cola, el símbolo mayor indica que la región crítica reside por completo en la cola derecha de la distribución de nuestro estadístico de prueba \bar{X} . ▮

Ejemplo 10.2: Un agente de bienes raíces afirma que 60% de todas las viviendas privadas que se construyen actualmente son casas con tres dormitorios. Para probar esta afirmación se inspecciona una muestra grande de viviendas nuevas. Se registra la proporción de las casas con 3 dormitorios y se utiliza como estadístico de prueba. Plantee las hipótesis nula y alternativa que se utilizarán en esta prueba y determine la ubicación de la región crítica.

Solución: Si el estadístico de prueba fuera considerablemente mayor o menor que $p = 0.6$, rechazaríamos la afirmación del agente. En consecuencia, deberíamos plantear las siguientes hipótesis:

$$H_0: p = 0.6,$$

$$H_1: p \neq 0.6.$$

La hipótesis alternativa implica una prueba de dos colas con la región crítica dividida por igual en ambas colas de la distribución de \hat{P} , nuestro estadístico de prueba. ▮

10.3 Uso de valores P para la toma de decisiones en la prueba de hipótesis

Al probar hipótesis en las que el estadístico de prueba es discreto, la región crítica se podría elegir de manera arbitraria y determinar su tamaño. Si α es demasiado grande, se reduce haciendo un ajuste en el valor crítico. Quizá sea necesario aumentar el tamaño

de la muestra para compensar la disminución que ocurre de manera automática en la potencia de la prueba.

Por generaciones enteras de análisis estadístico se ha vuelto costumbre elegir un α de 0.05 o 0.01 y seleccionar la región crítica de acuerdo con esto. Entonces, desde luego, el rechazo o no rechazo estrictos de H_0 dependerá de esa región crítica. Por ejemplo, si la prueba es de dos colas, α se fija a un nivel de significancia de 0.05 y el estadístico de prueba implica, digamos, la distribución normal estándar, entonces se observa un valor z de los datos y la región crítica es

$$z > 1.96 \quad \text{o} \quad z < -1.96,$$

donde el valor 1.96 corresponde a $z_{0.025}$ en la tabla A.3. Un valor de z en la región crítica sugiere la aseveración: "El valor del estadístico de prueba es significativo", el cual se puede traducir al lenguaje del caso. Por ejemplo, si la hipótesis es dada por

$$H_0: \mu = 10,$$

$$H_1: \mu \neq 10,$$

se puede decir: "La media difiere de manera significativa del valor 10".

Preselección de un nivel de significancia

Esta preselección de un nivel de significancia α tiene sus raíces en la filosofía de que se debe controlar el riesgo máximo de cometer un error tipo I. Sin embargo, este enfoque no explica los valores del estadístico de prueba que están "cerca" a la región crítica. Suponga, por ejemplo, que en el caso de $H_0: \mu = 10$, contra $H_1: \mu \neq 10$, se observa un valor $z = 1.87$. En términos estrictos, con $\alpha = 0.05$ el valor no es significativo; pero el riesgo de cometer un error tipo I si se rechaza H_0 en este caso difícilmente se podría considerar grave. De hecho, en una situación de dos colas, el riesgo se cuantifica como

$$P = 2P(Z > 1.87 \text{ cuando } \mu = 10) = 2(0.0307) = 0.0614.$$

Como resultado, 0.0614 es la probabilidad de obtener un valor de z tan grande o mayor (en magnitud) que 1.87 cuando, de hecho, $\mu = 10$. Aunque esta evidencia en contra de H_0 no es tan firme como la que resultaría de un rechazo a un nivel $\alpha = 0.05$, se trata de información importante para el usuario. De hecho, el uso continuo de $\alpha = 0.05$ o 0.01 tan sólo es un resultado de lo que los estándares han transmitido por generaciones. **En la estadística aplicada los usuarios han adoptado de forma extensa el método del valor P .** El método está diseñado para dar al usuario una alternativa (en términos de una probabilidad) a la mera conclusión de "rechazo" o "no rechazo". El cálculo del valor P también proporciona al usuario información importante cuando el valor z cae *dentro de la región crítica ordinaria*. Por ejemplo, si z es 2.73, resulta informativo para el usuario observar que

$$P = 2(0.0032) = 0.0064,$$

y, por consiguiente, el valor z es significativo a un nivel considerablemente menor que 0.05. Es importante saber que bajo la condición de H_0 un valor de $z = 2.73$ es un evento demasiado raro. A saber, un valor al menos tan grande en magnitud sólo ocurriría 64 veces en 10,000 experimentos.

Demstración gráfica de un valor P

Una manera muy simple de explicar gráficamente un valor P consiste en considerar dos muestras distintas. Suponga que se están considerando dos materiales para cubrir un tipo específico de metal con el fin de evitar la corrosión. Se obtienen especímenes y se cubre un grupo con el material 1 y otro grupo con el material 2. Los tamaños muestrales son $n_1 = n_2 = 10$ para cada muestra y la corrosión se mide en el porcentaje del área superficial afectada. La hipótesis plantea que las muestras provienen de distribuciones comunes con media $\mu = 10$. Supongamos que la varianza de la población es 1.0. Entonces, probamos

$$H_0: \mu_1 = \mu_2 = 10.$$

Representemos con la figura 10.8 una gráfica de puntos de los datos. Los datos se colocan en la distribución determinada por la hipótesis nula. Supongamos que los datos “x” se refieren al material 1 y que los datos “o” se refieren al material 2. Parece evidente que los datos realmente refutan la hipótesis nula. Pero, ¿cómo se podría resumir esto en un número? **El valor P se puede considerar simplemente como la probabilidad de obtener este conjunto de datos dado que las muestras provienen de la misma distribución.** Es evidente que esta probabilidad es muy pequeña, ¡digamos 0.00000001! Por consiguiente, el pequeño valor P evidentemente refuta H_0 , y la conclusión es que las medias de la población son significativamente diferentes.

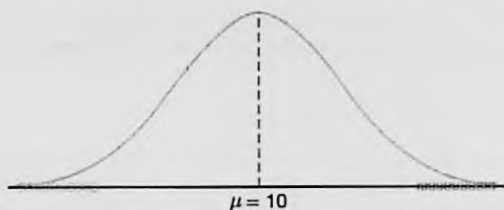


Figura 10.8: Datos que son probablemente generados de poblaciones que tienen dos medias diferentes.

El uso del método del valor P como auxiliar en la toma de decisiones es muy natural y casi todos los programas de cómputo que proporcionan el cálculo de pruebas de hipótesis ofrecen valores P , junto con valores del estadístico de prueba adecuado. La siguiente es una definición formal de un valor P .

Definición 10.5: Un valor P es el nivel (de significancia) más bajo en el que el valor observado del estadístico de prueba es significativo.

¿En qué difiere el uso de los valores P de la prueba de hipótesis clásica?

En este momento resulta tentador resumir los procedimientos que se asocian con la prueba de, digamos, $H_0: \theta = \theta_0$. Sin embargo, el estudiante que es novato en esta área deberá tener en cuenta que hay diferencias entre el enfoque y la filosofía del método

clásico de α fija, que tiene su momento más importante en la conclusión de “rechazar H_0 ” o “no rechazar H_0 ” y el método del valor P . En este último no se determina una α fija y las conclusiones se obtienen con base en el tamaño del valor P , según la apreciación subjetiva del ingeniero o del científico. Aun cuando los modernos programas de cómputo proporcionan valores P , es importante que el lector comprenda ambos enfoques para apreciar la totalidad de los conceptos. Por lo tanto, ofrecemos una breve lista con los pasos del procedimiento tanto para el método clásico como para el del valor P .

Aproximación a la prueba de hipótesis con probabilidad fija del error tipo I	<ol style="list-style-type: none"> 1. Establezca las hipótesis nula y alternativa. 2. Elija un nivel de significancia α fijo. 3. Seleccione un estadístico de prueba adecuado y establezca la región crítica con base en α. 4. Rechace H_0 si el estadístico de prueba calculado está en la región crítica. De otra manera, no rechace H_0. 5. Saque conclusiones científicas y de ingeniería.
Prueba de significancia (método del valor P)	<ol style="list-style-type: none"> 1. Establezca las hipótesis nula y alternativa. 2. Elija un estadístico de prueba adecuado. 3. Calcule el valor P con base en los valores calculados del estadístico de prueba. 4. Saque conclusiones con base en el valor P y los conocimientos del sistema científico.

En secciones posteriores de este capítulo y en los capítulos siguientes muchos ejemplos y ejercicios destacarán el método del valor P para obtener conclusiones científicas.

Ejercicios

10.1 Suponga que un alergólogo desea probar la hipótesis de que al menos 30% del público es alérgico a algunos productos de queso. Explique cómo el alergólogo podría cometer

- a) un error tipo I;
- b) un error tipo II.

10.2 Una socióloga se interesa en la eficacia de un curso de entrenamiento diseñado para lograr que más conductores utilicen los cinturones de seguridad en los automóviles.

- a) ¿Qué hipótesis pone a prueba si comete un error tipo I al concluir de manera errónea que el curso de entrenamiento no es eficaz?
- b) ¿Qué hipótesis pone a prueba si comete un error tipo II al concluir de forma errónea que el curso de entrenamiento es eficaz?

10.3 Se acusa a una empresa grande de discriminación en sus prácticas de contratación.

- a) ¿Qué hipótesis se pone a prueba si un jurado comete un error tipo I al encontrar culpable a la empresa?
- b) ¿Qué hipótesis se pone a prueba si un jurado comete un error tipo II al encontrar culpable a la empresa?

10.4 Un fabricante de telas considera que la proporción de pedidos de materia prima que llegan con retraso es $p = 0.6$. Si una muestra aleatoria de 10 pedidos indica que 3 o menos llegaron con retraso, la hipótesis de que $p = 0.6$ se debería rechazar a favor de la alternativa $p < 0.6$. Utilice la distribución binomial.

- a) Calcule la probabilidad de cometer un error tipo I si la proporción verdadera es $p = 0.6$.
- b) Calcule la probabilidad de cometer un error tipo II para las alternativas $p = 0.3$, $p = 0.4$ y $p = 0.5$.

10.5 Repita el ejercicio 10.4 pero suponga que se seleccionan 50 pedidos y que se define a la región crítica como $x \leq 24$, donde x es el número de pedidos en la muestra que llegaron con retraso. Utilice la aproximación normal.

10.6 Se estima que la proporción de adultos que vive en una pequeña ciudad que son graduados universitarios es $p = 0.6$. Para probar esta hipótesis se selecciona una muestra aleatoria de 15 adultos. Si el número de graduados en la muestra es cualquier número entre 6 y 12, no rechazaremos la hipótesis nula de que $p = 0.6$; de otro modo, concluiremos que $p \neq 0.6$.

- a) Evalúe α suponiendo que $p = 0.6$. Utilice la distribución binomial.

- b) Evalúe β para las alternativas $p = 0.5$ y $p = 0.7$.
 c) ¿Es éste un buen procedimiento de prueba?

10.7 Repita el ejercicio 10.6 pero suponga que se seleccionan 200 adultos y que la región de no rechazo se define como $110 \leq x \leq 130$, donde x es el número de individuos graduados universitarios en la muestra. Utilice la aproximación normal.

10.8 En la publicación *Relief from Arthritis* de Thorsons Publishers, Ltd., John E. Croft afirma que más de 40% de los individuos que sufren de osteoartritis experimentan un alivio medible con un ingrediente producido por una especie particular de mejillón que se encuentra en la costa de Nueva Zelanda. Para probar esa afirmación se suministra el extracto de mejillón a un grupo de 7 pacientes con osteoartritis. Si 3 o más de los pacientes experimentan alivio, no rechazaremos la hipótesis nula de que $p = 0.4$; de otro modo, concluiremos que $p < 0.4$.

- a) Evalúe α suponiendo que $p = 0.4$.
 b) Evalúe β para la alternativa $p = 0.3$.

10.9 Una tintorería afirma que un nuevo removedor de manchas quitará más de 70% de las manchas en las que se aplique. Para verificar esta afirmación el removedor de manchas se utilizará sobre 12 manchas elegidas al azar. Si se eliminan menos de 11 de las manchas, no se rechazará la hipótesis nula de que $p = 0.7$; de otra manera, concluiremos que $p > 0.7$.

- a) Evalúe α , suponiendo que $p = 0.7$.
 b) Evalúe β para la alternativa $p = 0.9$.

10.10 Repita el ejercicio 10.9 pero suponga que se tratan 100 manchas y que la región crítica se define como $x > 82$, donde x es el número de manchas eliminadas.

10.11 Repita el ejercicio 10.8 pero suponga que el extracto de mejillón se administra a 70 pacientes y que la región crítica se define como $x < 24$, donde x es el número de pacientes con osteoartritis que experimentan alivio.

10.12 Se pregunta a una muestra aleatoria de 400 votantes en cierta ciudad si están a favor de un impuesto adicional de 4% sobre las ventas de gasolina con el fin de obtener los fondos que se necesitan con urgencia para la reparación de calles. Si más de 220 votantes, pero menos de 260 de ellos, favorecen el impuesto sobre las ventas, concluiremos que 60% de los votantes lo apoyan.

- a) Calcule la probabilidad de cometer un error tipo I si 60% de los votantes están a favor del aumento de impuestos.
 b) ¿Cuál es la probabilidad de cometer un error tipo II al utilizar este procedimiento de prueba si en realidad sólo 48% de los votantes está a favor del impuesto adicional a la gasolina?

10.13 Suponga que en el ejercicio 10.12 concluimos que 60% de los votantes está a favor del impuesto sobre

las ventas de gasolina si más de 214 votantes, pero menos de 266 de ellos, lo favorecen. Demuestre que esta nueva región crítica tiene como resultado un valor más pequeño para α a costa de aumentar β .

10.14 Un fabricante desarrolla un nuevo sedal para pesca que, según afirma, tiene una resistencia media a la rotura de 15 kilogramos con una desviación estándar de 0.5 kilogramos. Para probar la hipótesis de que $\mu = 15$ kilogramos contra la alternativa de que $\mu < 15$ kilogramos se prueba una muestra aleatoria de 50 sedales. La región crítica se define como $\bar{x} < 14.9$.

- a) Calcule la probabilidad de cometer un error tipo I cuando H_0 es verdadera.
 b) Evalúe β para las alternativas $\mu = 14.8$ y $\mu = 14.9$ kilogramos.

10.15 En un restaurante de carnes una máquina de bebidas gaseosas se ajusta para que la cantidad de bebida que sirva se distribuya de forma aproximadamente normal, con una media de 200 mililitros y una desviación estándar de 15 mililitros. La máquina se verifica periódicamente tomando una muestra de 9 bebidas y calculando el contenido promedio. Si \bar{x} cae en el intervalo $191 < \bar{x} < 209$, se considera que la máquina opera de forma satisfactoria; de otro modo, se concluye que $\mu \neq 200$ mililitros.

- a) Calcule la probabilidad de cometer un error tipo I cuando $\mu = 200$ mililitros.
 b) Calcule la probabilidad de cometer un error tipo II cuando $\mu = 215$ mililitros.

10.16 Repita el ejercicio 10.15 para muestras de tamaño $n = 25$. Utilice la misma región crítica.

10.17 Se desarrolla un nuevo proceso de cura para cierto tipo de cemento que da como resultado una resistencia media a la compresión de 5000 kilogramos por centímetro cuadrado y una desviación estándar de 120 kilogramos. Para probar la hipótesis de que $\mu = 5000$ contra la alternativa de que $\mu < 5000$ se toma una muestra aleatoria de 50 piezas de cemento. La región crítica se define como $\bar{x} < 4970$.

- a) Calcule la probabilidad de cometer un error tipo I cuando H_0 es verdadera.
 b) Evalúe β para las alternativas $\mu = 4970$ y $\mu = 4960$.

10.18 Si graficamos las probabilidades de no rechazar H_0 que corresponden a diversas alternativas para μ (incluido el valor especificado para H_0) y conectamos todos los puntos mediante una curva suave, obtenemos la **curva característica de operación** del criterio de prueba o, simplemente, la curva CO. Observe que la probabilidad de no rechazar H_0 cuando es verdadera es simplemente $1 - \alpha$. Las curvas características de operación se utilizan con amplitud en aplicaciones industriales para proporcionar una muestra visual de los

méritos del criterio de prueba. Remítase al ejercicio 10.15 y calcule las probabilidades de no rechazar H_0 para los siguientes 9 valores de μ y grafique la curva CO: 184, 188, 192, 196, 200, 204, 208, 212 y 216.

10.4 Una sola muestra: pruebas respecto a una sola media

En esta sección consideramos de manera formal pruebas de hipótesis para una sola media de la población. Muchos de los ejemplos de las secciones anteriores incluyen pruebas sobre la media, por lo que el lector ya debería tener una idea de algunos de los detalles que aquí se describen.

Pruebas para una sola media (varianza conocida)

Primero deberíamos describir las suposiciones en las que se basa el experimento. El modelo para la situación subyacente se centra alrededor de un experimento con X_1, X_2, \dots, X_n , que representan una muestra aleatoria de una distribución con media μ y varianza $\sigma^2 > 0$. Considere primero la hipótesis

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

El estadístico de prueba adecuado se debe basar en la variable aleatoria \bar{X} . En el capítulo 8 se presentó el teorema del límite central, el cual establece en esencia que, sin importar la distribución de X , la variable aleatoria \bar{X} tiene una distribución casi normal con media μ y varianza σ^2/n para muestras de tamaño razonablemente grande. Por consiguiente, $\mu_{\bar{X}} = \mu$ y $\sigma_{\bar{X}}^2 = \sigma^2/n$. Podemos determinar, entonces, una región crítica basada en el promedio muestral calculado \bar{x} . Ahora ya debería quedarle claro al lector que habrá una región crítica de dos colas para la prueba.

Estandarización de \bar{X}

Es conveniente estandarizar \bar{X} e incluir de manera formal la variable aleatoria normal estándar Z , donde

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

Sabemos que, bajo H_0 , es decir, si $\mu = \mu_0$, entonces $\sqrt{n}(\bar{X} - \mu_0)/\sigma$ tiene una distribución $n(x; 0, 1)$ y, por lo tanto, la expresión

$$P\left(-z_{\alpha/2} < \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} < z_{\alpha/2}\right) = 1 - \alpha$$

se puede utilizar para escribir una región de no rechazo adecuada. El lector debería tener en la mente que, formalmente, la región crítica se diseña para controlar α , la probabilidad de cometer un error tipo I. Debería ser evidente que se necesita una *señal de evidencia de dos colas* para apoyar H_1 . Así, dado un valor calculado \bar{x} , la prueba formal implica rechazar H_0 si el *estadístico de prueba* z calculado cae en la región crítica que se describe a continuación.

Procedimiento de prueba para una sola media

(varianza conocida)

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} > z_{\alpha/2} \quad \text{o} \quad z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} < -z_{\alpha/2}$$

Si $-z_{\alpha/2} < z < z_{\alpha/2}$, no se rechaza H_0 . El rechazo de H_0 , desde luego, implica la aceptación de la hipótesis alternativa $\mu \neq \mu_0$. Con esta definición de la región crítica debería quedar claro que habrá α probabilidades de rechazar H_0 (al caer en la región crítica) cuando, en realidad, $\mu = \mu_0$.

Aunque es más fácil entender la región crítica escrita en términos de z , escribimos la misma región crítica en términos del promedio calculado \bar{x} . Lo siguiente se puede escribir como un procedimiento de decisión idéntico:

rechazar H_0 si $\bar{x} < a$ o $\bar{x} > b$,

donde

$$a = \mu_0 - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \quad b = \mu_0 + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

En consecuencia, para un nivel de significancia α , los valores críticos de la variable aleatoria z y \bar{x} se presentan en la figura 10.9.

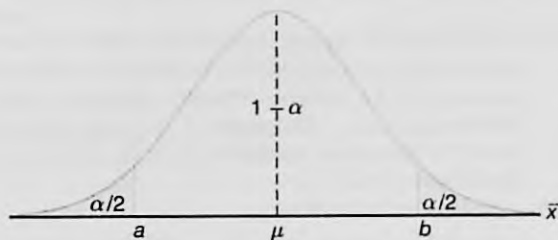


Figura 10.9: Región crítica para la hipótesis alternativa $\mu \neq \mu_0$.

Las pruebas de hipótesis unilaterales sobre la media incluyen el mismo estadístico que se describe en el caso bilateral. La diferencia, por supuesto, es que la región crítica sólo está en una cola de la distribución normal estándar. Por ejemplo, supongamos que buscamos probar

$$H_0: \mu = \mu_0,$$

$$H_1: \mu > \mu_0.$$

La señal que favorece H_1 proviene de valores grandes de z . Así, el rechazo de H_0 resulta cuando se calcula $z > z_\alpha$. Evidentemente, si la alternativa es $H_1: \mu < \mu_0$, la región crítica está por completo en la cola inferior, por lo que el rechazo resulta de $z < -z_\alpha$. Aunque en el caso de una prueba unilateral la hipótesis nula se puede escribir como $H_0: \mu \leq \mu_0$ o $H_0: \mu \geq \mu_0$, por lo general se escribe como $H_0: \mu = \mu_0$.

Los siguientes dos ejemplos ilustran pruebas de medias para el caso en el que se conoce σ .

Ejemplo 10.3: Una muestra aleatoria de 100 muertes registradas en Estados Unidos el año pasado reveló una vida promedio de 71.8 años. Si se supone una desviación estándar de la población de 8.9 años, ¿esto parece indicar que la vida media actual es mayor que 70 años? Utilice un nivel de significancia de 0.05.

Solución: 1. $H_0: \mu = 70$ años.

2. $H_1: \mu > 70$ años.

3. $\alpha = 0.05$.

4. Región crítica: $z > 1.645$, donde $z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$.

5. Cálculos: $\bar{x} = 71.8$ años, $\sigma = 8.9$ años, en consecuencia, $z = \frac{71.8 - 70}{8.9 / \sqrt{100}} = 2.02$.

6. Decisión: rechazar H_0 y concluir que la vida media actual es mayor que 70 años.

El valor P que corresponde a $z = 2.02$ es dado por el área de la región sombreada en la figura 10.10.

Si usamos la tabla A.3, tenemos

$$P = P(Z > 2.02) = 0.0217.$$

Como resultado, la evidencia a favor de H_1 es incluso más firme que la sugerida por un nivel de significancia de 0.05. J

Ejemplo 10.4: Un fabricante de equipo deportivo desarrolló un nuevo sedal para pesca sintético que, según afirma, tiene una resistencia media a la rotura de 8 kilogramos con una desviación estándar de 0.5 kilogramos. Pruebe la hipótesis de que $\mu = 8$ kilogramos contra la alternativa de que $\mu \neq 8$ kilogramos si se prueba una muestra aleatoria de 50 sedales y se encuentra que tienen una resistencia media a la rotura de 7.8 kilogramos. Utilice un nivel de significancia de 0.01.

Solución: 1. $H_0: \mu = 8$ kilogramos.

2. $H_1: \mu \neq 8$ kilogramos.

3. $\alpha = 0.01$.

4. Región crítica: $z < -2.575$ y $z > 2.575$, donde $z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$.

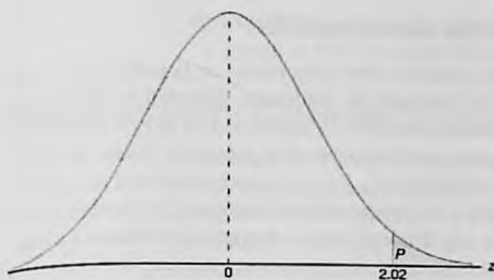
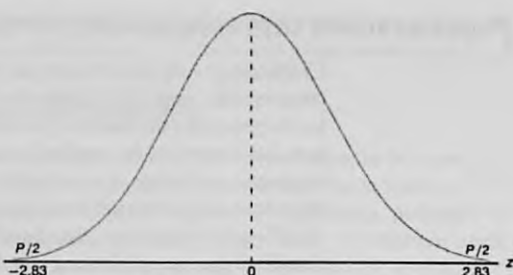
5. Cálculos: $\bar{x} = 7.8$ kilogramos, $n = 50$, en consecuencia, $z = \frac{7.8 - 8}{0.5 / \sqrt{50}} = -2.83$.

6. Decisión: rechazar H_0 y concluir que la resistencia promedio a la rotura no es igual a 8 sino que, de hecho, es menor que 8 kilogramos.

Como la prueba en este ejemplo es de dos colas, el valor de P que se desea es el doble del área de la región sombreada en la figura 10.11 a la izquierda de $z = -2.83$. Por lo tanto, si usamos la tabla A.3, tenemos

$$P = P(|Z| > 2.83) = 2P(Z < -2.83) = 0.0046,$$

que nos permite rechazar la hipótesis nula de que $\mu = 8$ kilogramos a un nivel de significancia menor que 0.01. J

Figura 10.10: Valor P para el ejemplo 10.3.Figura 10.11: Valor P para el ejemplo 10.4.

Relación con la estimación del intervalo de confianza

El lector ya se habrá dado cuenta de que el método de la prueba de hipótesis para la inferencia estadística de este capítulo está muy relacionado con el método del intervalo de confianza del capítulo 9. La estimación del intervalo de confianza incluye el cálculo de límites dentro de los cuales es "razonable" que resida el parámetro en cuestión. Para el caso de una sola media de la población μ con σ^2 conocida, la estructura tanto de la prueba de hipótesis como de la estimación del intervalo de confianza se basa en la variable aleatoria

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}.$$

Resulta que la prueba de $H_0: \mu = \mu_0$ contra $H_1: \mu \neq \mu_0$ a un nivel de significancia α es equivalente a calcular un intervalo de confianza del $100(1 - \alpha)\%$ sobre μ y rechazar H_0 , si μ_0 está fuera del intervalo de confianza. Si μ_0 está dentro del intervalo de confianza, no se rechaza la hipótesis. La equivalencia es muy intuitiva y se puede ilustrar de manera muy simple. Recuerde que con un valor observado \bar{x} , no rechazar H_0 a un nivel de significancia α implica que

$$-z_{\alpha/2} \leq \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \leq z_{\alpha/2},$$

que es equivalente a

$$\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu_0 \leq \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

La equivalencia de la estimación del intervalo de confianza con la prueba de hipótesis se extiende a las diferencias entre dos medias, varianzas, cocientes de varianzas, etcétera. Como resultado, el estudiante de estadística no debería considerar la estimación del intervalo de confianza y la prueba de hipótesis como formas separadas de inferencia estadística. Considere el ejemplo 9.2 de la página 271. El intervalo de confianza del 95% sobre la media es dado por los límites (2.50, 2.70). Por consiguiente, con la misma información muestral, no se rechazará una hipótesis bilateral sobre μ que incluya cualquier valor hipotético entre 2.50 y 2.70. A medida que exploremos diferentes áreas de la prueba de hipótesis seguiremos aplicando la equivalencia a la estimación del intervalo de confianza.

Pruebas sobre una sola media (varianza desconocida)

Ciertamente sospecharíamos que las pruebas sobre una media de la población μ con σ^2 desconocida, como la estimación del intervalo de confianza, deberían incluir el uso de la distribución t de Student. En términos estrictos, la aplicación de la t de Student tanto para los intervalos de confianza como para la prueba de hipótesis se desarrolla bajo los siguientes supuestos. Las variables aleatorias X_1, X_2, \dots, X_n representan una muestra aleatoria de una distribución normal con μ y σ^2 desconocidas. Entonces, la variable aleatoria $\sqrt{n}(\bar{X} - \mu)/S$ tiene una distribución t de Student con $n - 1$ grados de libertad. La estructura de la prueba es idéntica a la del caso en el que se conoce σ , excepto que el valor σ en el estadístico de prueba se reemplaza con el estimado calculado de S y la distribución normal estándar se reemplaza con una distribución t .

El estadístico t para la hipótesis bilateral para una prueba sobre una sola media (varianza desconocida)

$$H_0: \mu = \mu_0,$$

$$H_1: \mu \neq \mu_0,$$

rechazamos H_0 a un nivel de significancia α cuando el estadístico t calculado

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

excede a $t_{\alpha/2, n-1}$ o es menor que $-t_{\alpha/2, n-1}$.

El lector debería recordar de los capítulos 8 y 9 que la distribución t es simétrica alrededor del valor cero. Así, esta región crítica de dos colas se aplica de manera similar a la del caso en que se conoce σ . Para la hipótesis bilateral a un nivel de significancia α se aplican las regiones críticas de dos colas. Para $H_1: \mu > \mu_0$, el rechazo resulta cuando $t > t_{\alpha, n-1}$. Para $H_1: \mu < \mu_0$ la región crítica es dada por $T < -t_{\alpha, n-1}$.

Ejemplo 10.5: El Edison Electric Institute publica cifras del número de kilowatts-hora que gastan anualmente varios aparatos electrodomésticos. Se afirma que una aspiradora gasta un promedio de 46 kilowatts-hora al año. Si una muestra aleatoria de 12 hogares, que se incluye en un estudio planeado, indica que las aspiradoras gastan un promedio de 42 kilowatts-hora al año con una desviación estándar de 11.9 kilowatts-hora, ¿esto sugiere que las aspiradoras gastan, en promedio, menos de 46 kilowatts-hora al año a un nivel de significancia de 0.05? Suponga que la población de kilowatts-hora es normal.

Solución: 1. $H_0: \mu = 46$ kilowatts-hora.

2. $H_1: \mu < 46$ kilowatts-hora.

3. $\alpha = 0.05$.

4. Región crítica: $t < -1.796$, donde $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$ con 11 grados de libertad.

5. Cálculos: $\bar{x} = 42$ kilowatts-hora, $s = 11.9$ kilowatts-hora y $n = 12$.
En consecuencia,

$$t = \frac{42 - 46}{11.9/\sqrt{12}} = -1.16, \quad P = P(T < -1.16) \approx 0.135.$$

6. Decisión: no rechazar H_0 y concluir que el número promedio de kilowatts-hora que gastan al año las aspiradoras domésticas no es significativamente menor que 46. ■

Comentario sobre la prueba t de una sola muestra

Es probable que el lector haya observado que se mantiene la equivalencia de la prueba t de dos colas para una sola media y el cálculo de un intervalo de confianza sobre μ con σ reemplazada por s . Considere el ejemplo 9.5 de la página 275. En esencia, podemos ver ese cálculo como uno en el que encontramos todos los valores de μ_0 , el volumen medio hipotético de contenedores de ácido sulfúrico para los que la hipótesis $H_0: \mu = \mu_0$ no se rechazará con $\alpha = 0.05$. De nuevo, esto es consistente con el planteamiento: "si nos basamos en la información muestral, son razonables los valores del volumen medio de la población entre 9.74 y 10.26 litros".

En este punto vale la pena destacar algunos comentarios respecto a la suposición de normalidad. Indicamos que cuando se conoce σ , el teorema del límite central permite utilizar un estadístico de prueba o un intervalo de confianza que se base en Z , la variable aleatoria normal estándar. En términos estrictos, por supuesto, el teorema del límite central y, por lo tanto, el uso de la distribución normal estándar, no se aplica a menos que se conozca σ . En el capítulo 8 se estudió el desarrollo de la distribución t . Ahí se estableció que la normalidad sobre X_1, X_2, \dots, X_n era una suposición subyacente. Entonces, *en sentido estricto*, no se deberían utilizar las tablas de t de Student de puntos porcentuales para pruebas o intervalos de confianza, a menos que se sepa que la muestra proviene de una población normal. En la práctica rara vez se puede suponer un σ conocida. Sin embargo, se dispondría de una buena estimación a partir de experimentos anteriores. Muchos libros de estadística sugieren que, cuando $n \geq 30$, es posible reemplazar con seguridad σ por s en el estadístico de prueba

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

con una población que tiene forma de campana y aun así utilizar las tablas Z para la región crítica adecuada. Aquí la implicación es que en realidad se recurre al teorema del límite central y que se confía en el hecho de que $s \approx \sigma$. Evidentemente, cuando se hace esto el resultado debe considerarse como una aproximación. Por consiguiente, un valor P calculado (de la distribución Z) de 0.15 puede ser 0.12 o quizá 0.17; o un intervalo de confianza calculado puede ser un intervalo de 93% de confianza en vez de un intervalo de 95% como se desea. Entonces, ¿qué sucede en las situaciones donde $n \leq 30$? El usuario no puede confiar en que s se acerque a σ , y para tomar en cuenta la inexactitud de la estimación el intervalo de confianza debería ser más ancho o el valor crítico de mayor magnitud. Los puntos porcentuales de la distribución t logran esto, pero sólo son correctos cuando la muestra proviene de una distribución normal. Desde luego, se pueden utilizar las gráficas de probabilidad normal para tener cierta idea de la desviación de la normalidad en un conjunto de datos.

Para muestras pequeñas a menudo resulta difícil detectar desviaciones de una distribución normal. (Las pruebas de la bondad del ajuste se presentan en una sección posterior de este capítulo). Para distribuciones en forma de campana de las variables aleatorias X_1, X_2, \dots, X_n , es probable que el uso de la distribución t para pruebas o intervalos de confianza produzca resultados muy buenos. Cuando haya duda, el usuario debería recurrir a los procedimientos no paramétricos que se presentan en el capítulo 16.

Impresiones o salidas por computadora con comentarios para pruebas t de una sola muestra

Seguramente al lector le interesará ver comentarios impresos por computadora que muestren el resultado de una prueba t con una sola muestra. Suponga que un ingeniero se interesa en probar el sesgo en un medidor de pH. Se reúnen datos de una sustancia neutra ($\text{pH} = 7.0$). Se toma una muestra de las mediciones y los datos son los siguientes:

7.07 7.00 7.10 6.97 7.00 7.03 7.01 7.01 6.98 7.08

Entonces, es de interés probar

$$H_0: \mu = 7.0,$$

$$H_1: \mu \neq 7.0.$$

En este caso utilizamos el paquete de cómputo *MINITAB* para ilustrar el análisis del conjunto de datos anterior. Observe los componentes clave de la impresión o salida que se muestra en la figura 10.12. Desde luego, la media \bar{y} es 7.0250, StDev es simplemente la desviación estándar de la muestra $s = 0.044$ y SE Mean es el error estándar estimado de la media, y se calcula como $s/\sqrt{n} = 0.0139$. El valor t es el cociente

$$(7.0250 - 7) / 0.0139 = 1.80.$$

pH-meter									
7.07	7.00	7.10	6.97	7.00	7.03	7.01	7.01	6.98	7.08
MTB > Onet 'pH-meter'; SUBC> Test7.									
One-Sample T: pH-meter Test of mu = 7 vs not = 7									
Variable	N	Mean	StDev	SE Mean	95% CI	T	P		
pH-meter	10	7.02500	0.04403	0.01392	(6.99350, 7.05650)	1.80	0.106		

Figura 10.12: Impresión de *MINITAB* para la prueba t de una muestra para el medidor de pH.

El valor P de 0.106 sugiere resultados que no son concluyentes. No hay evidencia que sugiera un firme rechazo de H_0 (con base en una α de 0.05 o de 0.10), **ni se puede concluir con certeza que el medidor de pH esté libre de sesgo**. Observe que el tamaño de la muestra de 10 es muy pequeño. Un incremento en el tamaño de la muestra (quizás otro experimento) podría resolver las cosas. En la sección 10.6 aparece un análisis respecto al tamaño adecuado de la muestra.

10.5 Dos muestras: pruebas sobre dos medias

El lector deberá comprender la relación entre pruebas e intervalos de confianza y sólo puede confiar plenamente en los detalles que ofrece el material sobre el intervalo de confianza del capítulo 9. Las pruebas respecto a dos medias representan un conjunto de he-

ramientas analíticas muy importantes para el científico o el ingeniero. El procedimiento experimental es muy parecido al que se describe en la sección 9.8. Se extraen dos muestras aleatorias independientes de tamaños n_1 y n_2 , respectivamente, de dos poblaciones con medias μ_1 y μ_2 , y varianzas σ_1^2 y σ_2^2 . Sabemos que la variable aleatoria

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$$

tiene una distribución normal estándar. Suponemos aquí que n_1 y n_2 son suficientemente grandes, por lo que se aplica el teorema del límite central. Por supuesto, si las dos poblaciones son normales, el estadístico anterior tiene una distribución normal estándar incluso para n_1 y n_2 pequeñas. Evidentemente, si podemos suponer que $\sigma_1 = \sigma_2 = \sigma$, el estadístico anterior se reduce a

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sigma\sqrt{1/n_1 + 1/n_2}}.$$

Los dos estadísticos anteriores sirven como base para el desarrollo de los procedimientos de prueba que incluyen dos medias. La equivalencia entre las pruebas y los intervalos de confianza, junto con los detalles técnicos implicados en las pruebas sobre una media, permiten que la transición a pruebas con dos medias sea sencilla.

La hipótesis bilateral sobre dos medias se escribe de manera muy general como

$$H_0: \mu_1 - \mu_2 = d_0.$$

Es evidente que la alternativa puede ser bilateral o unilateral. De nuevo, la distribución que se utiliza es la distribución del estadístico de prueba bajo H_0 . Se calculan los valores \bar{x}_1 y \bar{x}_2 , y para σ_1 y σ_2 conocidas, el estadístico de prueba es dado por

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}},$$

con una región crítica de dos colas en el caso de una alternativa bilateral. Es decir, se rechaza H_0 a favor de $H_1: \mu_1 - \mu_2 \neq d_0$, si $z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$. Las regiones críticas de una cola se utilizan en el caso de alternativas unilaterales. El lector debería estudiar, como antes, el estadístico de prueba y estar satisfecho de que para, digamos $H_1: \mu_1 - \mu_2 > d_0$, la señal que favorece H_1 provenga de valores grandes de z . Por consiguiente, se aplica la región crítica de la cola superior.

Varianzas desconocidas pero iguales

Las situaciones más comunes que implican pruebas sobre dos medias son aquellas con varianzas desconocidas. Si el científico interesado está dispuesto a suponer que ambas distribuciones son normales y que $\sigma_1 = \sigma_2 = \sigma$, se puede utilizar la *prueba t agrupada* (a menudo llamada prueba *t* de dos muestras). El estadístico de prueba (véase la sección 9.8) es dado por el siguiente procedimiento de prueba.

Prueba t Para la hipótesis bilateral
agrupada de
dos muestras

$$H_0: \mu_1 = \mu_2,$$

$$H_1: \mu_1 \neq \mu_2,$$

rechazamos H_0 al nivel de significancia α cuando el estadístico t calculado

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{1/n_1 + 1/n_2}},$$

donde

$$s_p^2 = \frac{s_1^2(n_1 - 1) + s_2^2(n_2 - 1)}{n_1 + n_2 - 2}$$

excede a $t_{\alpha/2, n_1 + n_2 - 2}$ o es menor que $-t_{\alpha/2, n_1 + n_2 - 2}$.

Recuerde que en el capítulo 9 se explicó que los grados de libertad para la distribución t son un resultado del agrupamiento de la información de las dos muestras para estimar σ^2 . Las alternativas unilaterales, como era de esperarse, sugieren regiones críticas unilaterales. Por ejemplo, para $H_1: \mu_1 - \mu_2 > d_0$, rechace $H_0: \mu_1 - \mu_2 = d_0$ cuando $t > t_{\alpha, n_1 + n_2 - 2}$.

Ejemplo 10.6: Se llevó a cabo un experimento para comparar el desgaste por abrasivos de dos diferentes materiales laminados. Se probaron 12 piezas del material 1 exponiendo cada pieza a una máquina para medir el desgaste. Se probaron 10 piezas del material 2 de manera similar. En cada caso se observó la profundidad del desgaste. Las muestras del material 1 revelaron un desgaste promedio (codificado) de 85 unidades con una desviación estándar muestral de 4; en tanto que las muestras del material 2 revelaron un promedio de 81 y una desviación estándar muestral de 5. ¿Podríamos concluir, a un nivel de significancia de 0.05, que el desgaste abrasivo del material 1 excede al del material 2 en más de 2 unidades? Suponga que las poblaciones son aproximadamente normales con varianzas iguales.

Solución: Representemos con μ_1 y μ_2 las medias de la población del desgaste abrasivo para el material 1 y el material 2, respectivamente.

1. $H_0: \mu_1 - \mu_2 = 2.$

2. $H_1: \mu_1 - \mu_2 > 2.$

3. $\alpha = 0.05.$

4. Región crítica: $t > 1.725$, donde $t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{1/n_1 + 1/n_2}}$ con $v = 20$ grados de libertad.

5. Cálculos:

$$\bar{x}_1 = 85, \quad s_1 = 4, \quad n_1 = 12,$$

$$\bar{x}_2 = 81, \quad s_2 = 5, \quad n_2 = 10.$$

En consecuencia,

$$s_p = \sqrt{\frac{(11)(16) + (9)(25)}{12 + 10 - 2}} = 4.478,$$

$$t = \frac{(85 - 81) - 2}{4.478\sqrt{1/12 + 1/10}} = 1.04,$$

$$P = P(T > 1.04) \approx 0.16. \quad (\text{Véase la tabla A.4}).$$

6. Decisión: no rechazar H_0 . No podemos concluir que el desgaste abrasivo del material 1 excede al del material 2 en más de 2 unidades. ■

Varianzas desconocidas pero diferentes

Hay situaciones donde al analista no le es posible suponer que $\sigma_1 = \sigma_2$. De la sección 9.8 recuerde que, si las poblaciones son normales, el estadístico

$$T' = \frac{(\bar{X}_1 - \bar{X}_2) - d_0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

tiene una distribución t aproximada con grados de libertad aproximados

$$v = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{(s_1^2/n_1)^2/(n_1 - 1) + (s_2^2/n_2)^2/(n_2 - 1)}.$$

Como resultado, el procedimiento de prueba consiste en *no rechazar* H_0 cuando

$$-t_{\alpha/2, v} < t' < t_{\alpha/2, v},$$

con v dado como antes. De nuevo, como en el caso de la prueba t agrupada, las alternativas unilaterales sugieren regiones críticas unilaterales.

Observaciones pareadas

Un estudio de la prueba t de dos muestras o el intervalo de confianza sobre la diferencia entre medias deberían sugerir la necesidad de un diseño experimental. Recuerde el análisis de las unidades experimentales en el capítulo 9, donde se sugirió que las condiciones de las dos poblaciones (a menudo denominadas como los dos tratamientos) se deberían asignar de manera aleatoria a las unidades experimentales. Esto se realiza para evitar resultados sesgados debido a las diferencias sistemáticas entre unidades experimentales. En otras palabras, en términos de la jerga para la prueba de hipótesis, es importante que la diferencia significativa que se encuentre entre las medias se deba a las diferentes condiciones de las poblaciones y no a las unidades experimentales en el estudio. Por ejemplo, considere el ejercicio 9.40 de la sección 9.9. Los 20 tallos desempeñan el papel de unidades experimentales. Diez de ellos se tratan con nitrógeno y 10 se dejan sin tratamiento. Es muy importante que esta asignación a los tratamientos "con nitrógeno" y "sin nitrógeno" sea aleatoria para garantizar que las diferencias sistemáticas entre los tallos no interfieran con una comparación válida entre las medias.

En el ejemplo 10.6 el momento de la medición es la opción más probable de la unidad experimental. Las 22 piezas de material se deberían medir en orden aleatorio.

Necesitamos protegernos contra la posibilidad de que las mediciones del desgaste que se realicen casi al mismo tiempo tiendan a dar resultados similares. *No se esperan diferencias sistemáticas* (no aleatorias) en las unidades experimentales. Sin embargo, las asignaciones aleatorias protegen contra el problema.

Las referencias a la planeación de experimentos, aleatorización, elección del tamaño de la muestra, etcétera, continuarán influyendo en gran parte del desarrollo en los capítulos 13, 14 y 15. Cualquier científico o ingeniero cuyo interés resida en el análisis de datos reales debería estudiar este material. La prueba t agrupada se amplía en el capítulo 13 para cubrir más de dos medias.

La prueba de dos medias se puede llevar a cabo cuando los datos están en forma de observaciones pareadas, como se estudió en el capítulo 9. En esta estructura de pareado las condiciones de las dos poblaciones (tratamientos) se asignan de forma aleatoria dentro de unidades homogéneas. El cálculo del intervalo de confianza para $\mu_1 - \mu_2$ en la situación con observaciones pareadas se basa en la variable aleatoria

$$T = \frac{\bar{D} - \mu_D}{S_d/\sqrt{n}},$$

donde \bar{D} y S_d son variables aleatorias que representan la media muestral y la desviación estándar de las diferencias de las observaciones en las unidades experimentales. Como en el caso de la *prueba t agrupada*, la suposición es que las observaciones de cada población son normales. Este problema de dos muestras se reduce en esencia a un problema de una muestra utilizando las diferencias calculadas d_1, d_2, \dots, d_n . Por consiguiente, la hipótesis se reduce a

$$H_0: \mu_D = d_0.$$

El estadístico de prueba calculado es dado entonces por

$$t = \frac{\bar{d} - d_0}{s_d/\sqrt{n}}.$$

Las regiones críticas se construyen usando la distribución t con $n - 1$ grados de libertad.

El problema de la interacción en una prueba t pareada

El siguiente estudio de caso no sólo ilustra el uso de la prueba t pareada, sino que el análisis revelará mucho sobre las dificultades que surgen cuando ocurre una interacción entre los tratamientos y las unidades experimentales en la estructura de la t pareada. Recuerde que la interacción entre factores se presentó en la sección 1.7, en un análisis de los tipos generales de estudios estadísticos. El concepto de interacción será un tema importante desde el capítulo 13 hasta el 15.

Existen ciertos tipos de pruebas estadísticas en los que la existencia de una interacción produce dificultades. Un ejemplo es la prueba t pareada. En la sección 9.9 se utilizó la estructura pareada en el cálculo de un intervalo de confianza sobre la diferencia entre dos medias, y se reveló la ventaja del pareado para situaciones en que las unidades experimentales son homogéneas. El pareado produce una reducción en σ_D , la desviación estándar de una diferencia $D_i = X_{1i} - X_{2i}$, como se explicó en la sección 9.9. Si hay una interacción entre los tratamientos y las unidades experimentales, la ventaja lograda

mediante el pareado se podría reducir de manera sustancial. Por consiguiente, en el ejemplo 9.13 de la página 293 la suposición de la ausencia de interacción permitió que la diferencia en los niveles medios de TCDD (plasma contra tejido adiposo) fuera la misma en todos los veteranos. Un vistazo rápido a los datos sugiere que no hay una violación significativa de los supuestos de ausencia de interacción.

Para demostrar cómo influye la interacción en $\text{Var}(D)$ y, por lo tanto, en la calidad de la prueba t pareada, es aleccionador revisar la i -ésima diferencia dada por $D_i = X_{1i} - X_{2i} = (\mu_1 - \mu_2) + (\epsilon_{1i} - \epsilon_{2i})$, donde X_{1i} y X_{2i} se toman de la i -ésima unidad experimental. Si la unidad pareada es homogénea, los errores en X_{1i} y en X_{2i} serán similares y no independientes. En el capítulo 9 señalamos que la covarianza positiva entre los errores da como resultado una $\text{Var}(D)$ reducida. Por consiguiente, el tamaño de la diferencia en los tratamientos y la relación entre los errores en X_{1i} y X_{2i} , a los que contribuye la unidad experimental, tenderán a permitir la detección de una diferencia significativa.

¿Qué condiciones resultan en una interacción?

Consideremos una situación en la que las unidades experimentales no son homogéneas. Más bien, considere la i -ésima unidad experimental con las variables aleatorias X_{1i} y X_{2i} que no son similares. Sean ϵ_{1i} y ϵ_{2i} variables aleatorias que representan los errores en los valores X_{1i} y X_{2i} , respectivamente, en la unidad i -ésima. Así, podemos escribir

$$X_{1i} = \mu_1 + \epsilon_{1i} \text{ y } X_{2i} = \mu_2 + \epsilon_{2i}.$$

Los errores con valor esperado cero podrían tender a provocar que los valores de respuesta X_{1i} y X_{2i} se muevan en direcciones opuestas, dando como resultado un valor negativo para $\text{Cov}(\epsilon_{1i}, \epsilon_{2i})$ y, por ende, un valor negativo para $\text{Cov}(X_{1i}, X_{2i})$. En realidad, el modelo se podría volver aún más complicado por el hecho de que $\sigma_1^2 = \text{Var}(\epsilon_{1i}) \neq \sigma_2^2 = \text{Var}(\epsilon_{2i})$. Los parámetros de la varianza y la covarianza podrían variar entre las n unidades experimentales. Así, a diferencia del caso con homogeneidad, D_i tenderá a ser muy diferente en todas las unidades experimentales debido a la naturaleza heterogénea de la diferencia en $\epsilon_1 - \epsilon_2$ entre las unidades. Esto produce la interacción entre los tratamientos y las unidades. Además, para una unidad experimental específica (véase el teorema 4.9),

$$\sigma_D^2 = \text{Var}(D) = \text{Var}(\epsilon_1) + \text{Var}(\epsilon_2) - 2 \text{Cov}(\epsilon_1, \epsilon_2)$$

está inflado por el término negativo de covarianza, de manera que la ventaja lograda por el pareado en el caso de la unidad homogénea se pierde en el caso que aquí se describe. En tanto que la inflación en $\text{Var}(D)$ variará de un caso a otro, en algunas situaciones existe el peligro de que el aumento en la varianza neutralice cualquier diferencia que exista entre μ_1 y μ_2 . Desde luego, un valor grande de \bar{d} en el estadístico t podría reflejar una diferencia en el tratamiento que compense el estimado inflado de la varianza s_d^2 .

Estudio de caso 10.1: **Datos de muestra de sangre:** En un estudio realizado en el Departamento de Silvicultura y Fauna de Virginia Tech, J. A. Wesson examinó la influencia del fármaco *succinylcholine* sobre los niveles de circulación de andrógenos en la sangre. Se obtuvieron muestras de sangre de venados salvajes inmediatamente después de recibir una inyección intramuscular de *succinylcholine* con dardos de un rifle de caza. Treinta minutos después se obtuvo una segunda muestra de sangre y después los venados fueron liberados. Los

niveles de andrógenos de 15 venados al momento de la captura y 30 minutos más tarde, medidos en nanogramos por mililitro (ng/mL), se presentan en la tabla 10.2.

Suponga que las poblaciones de niveles de andrógenos al momento de la inyección y 30 minutos después se distribuyen normalmente, y pruebe, a un nivel de significancia de 0.05, si las concentraciones de andrógenos se alteraron después de 30 minutos.

Tabla 10.2: Datos para el estudio de caso 10.1

Venado	Andrógenos (ng/mL)		d_i
	Al momento de la inyección	30 minutos después de la inyección	
1	2.76	7.02	4.26
2	5.18	3.10	-2.08
3	2.68	5.44	2.76
4	3.05	3.99	0.94
5	4.10	5.21	1.11
6	7.05	10.26	3.21
7	6.60	13.91	7.31
8	4.79	18.53	13.74
9	7.39	7.91	0.52
10	7.30	4.85	-2.45
11	11.78	11.10	-0.68
12	3.90	3.74	-0.16
13	26.00	94.03	68.03
14	67.48	94.03	26.55
15	17.04	41.70	24.66

Solución: Sean μ_1 y μ_2 la concentración promedio de andrógenos al momento de la inyección y 30 minutos después, respectivamente. Procedemos como sigue:

- $H_0: \mu_1 = \mu_2$ o $\mu_D = \mu_1 - \mu_2 = 0$.
- $H_1: \mu_1 \neq \mu_2$ o $\mu_D = \mu_1 - \mu_2 \neq 0$.
- $\alpha = 0.05$.
- Región crítica: $t < -2.145$ y $t > 2.145$, donde $t = \frac{\bar{d} - d_0}{s_D / \sqrt{n}}$ con $v = 14$ grados de libertad.
- Cálculos: La media muestral y la desviación estándar para las d_i son

$$\bar{d} = 9.848 \text{ y } s_d = 18.474.$$

Por lo tanto,

$$t = \frac{9.848 - 0}{18.474 / \sqrt{15}} = 2.06.$$

- Aunque el estadístico t no es significativo al nivel 0.05, de la tabla A.4,

$$P = P(|T| > 2.06) \approx 0.06.$$

Como resultado, existe cierta evidencia de que hay una diferencia en los niveles medios circulantes de andrógenos.

La suposición de la ausencia de interacción implicaría que el efecto sobre los niveles de andrógenos de los venados es casi el mismo en los datos de ambos tratamientos, es decir, en el momento de la inyección de *succinylcholine* y 30 minutos después. Esto se puede expresar cambiando los papeles de los dos factores; por ejemplo, la diferencia en los tratamientos es casi igual en todas las unidades, es decir, los venados. Ciertamente hay algunas combinaciones venado/tratamiento para las que parece ser válida la suposición de ausencia de interacción, pero difícilmente existen evidencias firmes de que las unidades experimentales sean homogéneas. Sin embargo, la naturaleza de la interacción y el incremento resultante en $\text{Var}(\bar{D})$ parecen estar dominados por una diferencia sustancial en los tratamientos. Esto también es demostrado por el hecho de que 11 de los 15 venados mostraron señales positivas para las d_i calculadas y las d_i negativas (para los venados 2, 10, 11 y 12) son pequeñas en magnitud comparadas con las 12 positivas. Por consiguiente, al parecer el nivel medio de andrógenos es significativamente más alto 30 minutos después de la inyección que en el momento en que se aplica, y las conclusiones podrían ser más firmes de lo que sugiere $p = 0.06$.

Impresiones por computadora con comentarios para pruebas t pareadas

La figura 10.13 presenta una impresión por computadora del SAS para una prueba t pareada usando los datos del estudio de caso 10.1. Observe que el listado se parece al de una prueba t de una sola muestra y, por supuesto, esto es con exactitud lo que se realizó, ya que la prueba busca determinar si \bar{d} es significativamente diferente de cero.

Analysis Variable : Diff				
N	Mean	Std Error	t Value	Pr > t
15	9.8480000	4.7698699	2.06	0.0580

Figura 10.13: Impresión por computadora del SAS de la prueba t pareada para los datos del estudio de caso 10.1.

Resumen de los procedimientos de prueba

Mientras completamos el desarrollo formal de pruebas sobre medias de la población, ofrecemos la tabla 10.3, que resume el procedimiento de prueba para los casos de una sola media y de dos medias. Observe el procedimiento aproximado cuando las distribuciones son normales y las varianzas se desconocen pero no se suponen iguales. Este estadístico se estudió en el capítulo 9.

10.6 Elección del tamaño de la muestra para la prueba de medias

En la sección 10.2 demostramos cómo el analista puede explotar las relaciones entre el tamaño de la muestra, el nivel de significancia α y la potencia de la prueba para alcanzar cierto estándar de calidad. En la mayoría de las circunstancias prácticas el experimento debería planearse y, de ser posible, elegir el tamaño de la muestra antes del proceso de recolección de datos. Por lo general el tamaño de la muestra se determina de modo que

Tabla 10.3: Pruebas relacionadas con medias

H_0	Valor del estadístico de prueba	H_1	Región crítica
$\mu = \mu_0$	$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$; σ conocida	$\mu < \mu_0$ $\mu > \mu_0$ $\mu \neq \mu_0$	$z < -z_\alpha$ $z > z_\alpha$ $z < -z_{\alpha/2}$ o $z > z_{\alpha/2}$
$\mu = \mu_0$	$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$; $v = n - 1$, σ desconocida	$\mu < \mu_0$ $\mu > \mu_0$ $\mu \neq \mu_0$	$t < -t_\alpha$ $t > t_\alpha$ $t < -t_{\alpha/2}$ o $t > t_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$; σ_1 y σ_2 conocidas	$\mu_1 - \mu_2 < d_0$ $\mu_1 - \mu_2 > d_0$ $\mu_1 - \mu_2 \neq d_0$	$z < -z_\alpha$ $z > z_\alpha$ $z < -z_{\alpha/2}$ o $z > z_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{1/n_1 + 1/n_2}}$; $v = n_1 + n_2 - 2$, $\sigma_1 = \sigma_2$ pero desconocidas $s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$	$\mu_1 - \mu_2 < d_0$ $\mu_1 - \mu_2 > d_0$ $\mu_1 - \mu_2 \neq d_0$	$t < -t_\alpha$ $t > t_\alpha$ $t < -t_{\alpha/2}$ o $t > t_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$t' = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$; $v = \frac{(s_1^2/n_1)^2}{n_1 - 1} + \frac{(s_2^2/n_2)^2}{n_2 - 1}$, $\sigma_1 \neq \sigma_2$ y desconocidas	$\mu_1 - \mu_2 < d_0$ $\mu_1 - \mu_2 > d_0$ $\mu_1 - \mu_2 \neq d_0$	$t' < -t_\alpha$ $t' > t_\alpha$ $t' < -t_{\alpha/2}$ o $t' > t_{\alpha/2}$
$\mu_D = d_0$ observaciones pareadas	$t = \frac{\bar{d} - d_0}{s_d/\sqrt{n}}$; $v = n - 1$	$\mu_D < d_0$ $\mu_D > d_0$ $\mu_D \neq d_0$	$t < -t_\alpha$ $t > t_\alpha$ $t < -t_{\alpha/2}$ o $t > t_{\alpha/2}$

permita lograr una buena potencia para una α fija y una alternativa específica fija. Esta alternativa fija puede estar en la forma de $\mu - \mu_0$ en el caso de una hipótesis que incluya una sola media o $\mu_1 - \mu_2$ en el caso de un problema que implique dos medias. Los casos específicos serán ilustrativos.

Suponga que deseamos probar la hipótesis

$$H_0: \mu = \mu_0,$$

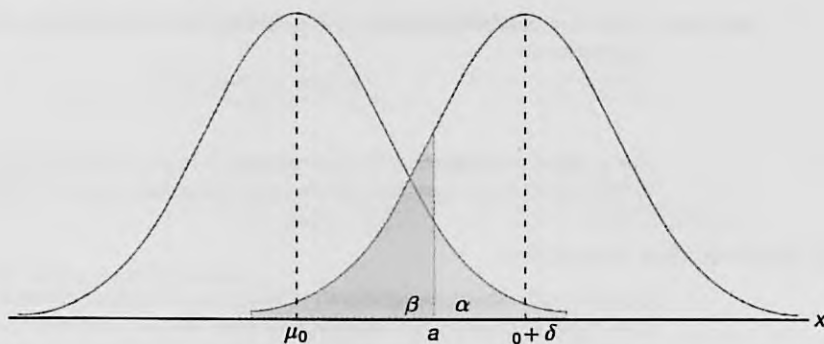
$$H_1: \mu > \mu_0,$$

con un nivel de significancia α , cuando se conoce la varianza σ^2 . Para una alternativa específica, digamos, $\mu = \mu_0 + \delta$, en la figura 10.14 se muestra que la potencia de nuestra prueba es

$$1 - \beta = P(\bar{X} > a \text{ cuando } \mu = \mu_0 + \delta).$$

Por lo tanto,

$$\begin{aligned} \beta &= P(\bar{X} < a \text{ cuando } \mu = \mu_0 + \delta) \\ &= P\left[\frac{\bar{X} - (\mu_0 + \delta)}{\sigma/\sqrt{n}} < \frac{a - (\mu_0 + \delta)}{\sigma/\sqrt{n}} \text{ cuando } \mu = \mu_0 + \delta\right]. \end{aligned}$$

Figura 10.14: Prueba de $\mu = \mu_0$ contra $\mu = \mu_0 + \delta$.

Bajo la hipótesis alternativa $\mu = \mu_0 + \delta$, el estadístico

$$\frac{\bar{X} - (\mu_0 + \delta)}{\sigma/\sqrt{n}}$$

es la variable normal estándar Z . Por lo tanto,

$$\beta = P\left(Z < \frac{a - \mu_0}{\sigma/\sqrt{n}} - \frac{\delta}{\sigma/\sqrt{n}}\right) = P\left(Z < z_\alpha - \frac{\delta}{\sigma/\sqrt{n}}\right),$$

de donde concluimos que

$$-z_\beta = z_\alpha - \frac{\delta\sqrt{n}}{\sigma},$$

y, en consecuencia,

$$\text{Elección del tamaño de la muestra: } n = \frac{(z_\alpha + z_\beta)^2 \sigma^2}{\delta^2},$$

un resultado que también es verdadero cuando la hipótesis alternativa es $\mu < \mu_0$.

En el caso de una prueba de dos colas obtenemos la potencia $1 - \beta$ para una alternativa específica cuando

$$n \approx \frac{(z_{\alpha/2} + z_\beta)^2 \sigma^2}{\delta^2}.$$

Ejemplo 10.7: Suponga que deseamos probar la hipótesis

$$H_0: \mu = 68 \text{ kilogramos,}$$

$$H_1: \mu > 68 \text{ kilogramos,}$$

para los pesos de estudiantes hombres en cierta universidad usando un nivel de significancia $\alpha = 0.05$ cuando se sabe que $\sigma = 5$. Calcule el tamaño muestral que se requiere si la potencia de nuestra prueba debe ser 0.95 cuando la media real es 69 kilogramos.

Solución: Como $\alpha = \beta = 0.05$, tenemos $z_\alpha = z_\beta = 1.645$. Para la alternativa $\beta = 69$ tomamos $\delta = 1$ y entonces,

$$n = \frac{(1.645 + 1.645)^2(25)}{1} = 270.6.$$

Por lo tanto, se requieren 271 observaciones si la prueba debe rechazar la hipótesis nula el 95% de las veces cuando, de hecho, μ es tan grande como 69 kilogramos. \blacksquare

El caso de dos muestras

Se puede utilizar un procedimiento similar para determinar el tamaño de la muestra $n = n_1 = n_2$ que se requiere para una potencia específica de la prueba en que se comparan dos medias de la población. Por ejemplo, suponga que deseamos probar la hipótesis

$$H_0: \mu_1 - \mu_2 = d_0,$$

$$H_1: \mu_1 - \mu_2 \neq d_0,$$

cuando se conocen σ_1 y σ_2 . Para una alternativa específica, digamos $\mu_1 - \mu_2 = d_0 + \delta$, en la figura 10.15 se muestra que la potencia de nuestra prueba es

$$1 - \beta = P(|\bar{X}_1 - \bar{X}_2| > a \text{ cuando } \mu_1 - \mu_2 = d_0 + \delta).$$

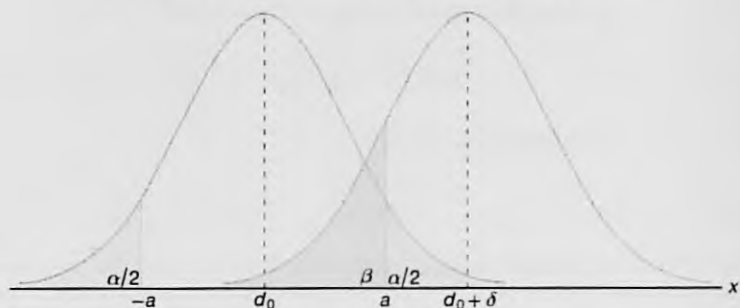


Figura 10.15: Prueba de $\mu_1 - \mu_2 = d_0$ contra $\mu_1 - \mu_2 = d_0 + \delta$.

Por lo tanto,

$$\begin{aligned} \beta &= P(-a < \bar{X}_1 - \bar{X}_2 < a \text{ cuando } \mu_1 - \mu_2 = d_0 + \delta) \\ &= P\left[\frac{-a - (d_0 + \delta)}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}} < \frac{(\bar{X}_1 - \bar{X}_2) - (d_0 + \delta)}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}} \right. \\ &\quad \left. < \frac{a - (d_0 + \delta)}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}} \text{ cuando } \mu_1 - \mu_2 = d_0 + \delta \right]. \end{aligned}$$

Con la hipótesis alternativa $\mu_1 - \mu_2 = d_0 + \delta$, el estadístico

$$\frac{\bar{X}_1 - \bar{X}_2 - (d_0 + \delta)}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}}$$

es la variable normal estándar Z . Ahora bien, al escribir

$$-z_{\alpha/2} = \frac{-a - d_0}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}} \quad \text{y} \quad z_{\alpha/2} = \frac{a - d_0}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}},$$

tenemos

$$\beta = P \left[-z_{\alpha/2} - \frac{\delta}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}} < Z < z_{\alpha/2} - \frac{\delta}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}} \right],$$

de donde concluimos que

$$-z_{\beta} \approx z_{\alpha/2} - \frac{\delta}{\sqrt{(\sigma_1^2 + \sigma_2^2)/n}},$$

y, por lo tanto,

$$n \approx \frac{(z_{\alpha/2} + z_{\beta})^2 (\sigma_1^2 + \sigma_2^2)}{\delta^2}.$$

Para la prueba de una sola cola, la expresión para el tamaño requerido de la muestra cuando $n = n_1 = n_2$ es

$$\text{Elección del tamaño de la muestra: } n = \frac{(z_{\alpha} + z_{\beta})^2 (\sigma_1^2 + \sigma_2^2)}{\delta^2}.$$

Cuando se desconoce la varianza de la población (o varianzas en la situación de dos muestras), la elección del tamaño de la muestra no es directa. Al probar la hipótesis $\mu = \mu_0$ cuando el valor verdadero es $\mu = \mu_0 + \delta$, el estadístico

$$\frac{\bar{X} - (\mu_0 + \delta)}{S/\sqrt{n}}$$

no sigue la distribución t , como se podría esperar, más bien sigue la **distribución t no central**. Sin embargo, existen tablas o gráficas que se basan en la distribución t no central para determinar el tamaño adecuado de la muestra, si se dispone de algún estimado de σ o si δ es un múltiplo de σ . La tabla A.8 proporciona los tamaños muestrales necesarios para controlar los valores de α y β para diversos valores de

$$\Delta = \frac{|\delta|}{\sigma} = \frac{|\mu - \mu_0|}{\sigma}$$

en el caso de pruebas de una y de dos colas. En el caso de la prueba t de dos muestras en la que se desconocen las varianzas pero se suponen iguales, obtenemos los tamaños muestrales $n = n_1 = n_2$ necesarios para controlar los valores de α y β para diversos valores de

$$\Delta = \frac{|\delta|}{\sigma} = \frac{|\mu_1 - \mu_2 - d_0|}{\sigma}$$

de la tabla A.9.

Ejemplo 10.8: Al comparar el comportamiento de dos catalizadores sobre el efecto del producto de una reacción se realiza una prueba t de dos muestras con $\alpha = 0.05$. Se considera que las

varianzas de los productos son iguales para los dos catalizadores. ¿De qué tamaño debe ser una muestra para cada catalizador si se desea probar la hipótesis

$$H_0: \mu_1 = \mu_2,$$

$$H_1: \mu_1 \neq \mu_2$$

si es esencial detectar una diferencia de 0.8σ entre los catalizadores con 0.9 de probabilidad?

Solución: De la tabla A.9, con $\alpha = 0.05$ para una prueba de dos colas, $\beta = 0.1$ y

$$\Delta = \frac{|0.8\sigma|}{\sigma} = 0.8,$$

encontramos que el tamaño requerido de la muestra es $n = 34$.

En situaciones prácticas sería difícil forzar a un científico o a un ingeniero a hacer un compromiso sobre la información a partir de la cual se puede encontrar un valor de Δ . Se recuerda al lector que el valor Δ cuantifica el tipo de diferencia entre las medias que el científico considera importantes; es decir, una diferencia que se considere *significativa* desde un punto de vista científico, no estadístico. El ejemplo 10.8 ilustra cómo suele hacerse esta elección, a saber, mediante la selección de una fracción de σ . Evidentemente, si el tamaño de la muestra se basa en una elección de $|\delta|$, que es una fracción pequeña de σ , el tamaño muestral que resulta podría ser muy grande comparado con lo que permite el estudio.

10.7 Métodos gráficos para comparar medias

En el capítulo 1 se puso mucha atención a la presentación de datos en forma gráfica, como los diagramas de tallo y hojas y las gráficas de caja y bigote. En la sección 8.8 las gráficas de cuantiles y las gráficas normales cuantil-cuantil se utilizaron para brindar una “imagen” y resumir así un conjunto de datos experimentales. Muchos paquetes de cómputo producen representaciones gráficas. A medida que procedamos con otras formas de análisis de datos, por ejemplo, el análisis de regresión y el análisis de varianza, los métodos gráficos se vuelven aún más informativos.

Los auxiliares gráficos no se pueden utilizar como un reemplazo del propio procedimiento de prueba. En realidad, el valor del estadístico de prueba indica el tipo adecuado de evidencia en apoyo de H_0 o H_1 . Sin embargo, una imagen ofrece una buena ilustración y a menudo es un mejor comunicador de evidencia para el beneficiario del análisis. Además, una imagen con frecuencia dejará claro por qué se encontró una diferencia significativa. La falla de una suposición importante se puede expresar mediante un resumen gráfico.

Para la comparación de medias, las gráficas de caja y bigote simultáneas proporcionan una imagen clara. El lector debería recordar que estas gráficas muestran el percentil 25, el percentil 75 y la mediana en un conjunto de datos. Además, las extensiones muestran los extremos en un conjunto de datos. Considere el ejercicio 10.40 al final de esta sección. Se midieron los niveles en plasma de ácido ascórbico en dos grupos de mujeres embarazadas: fumadoras y no fumadoras. En la figura 10.16 se observan las gráficas de caja y bigote para ambos grupos de mujeres y dos cosas son muy evidentes; al tomar en cuenta la variabilidad parece haber una diferencia despreciable en las medias muestrales. Además, parece que la variabilidad en los dos grupos es hasta cierto punto diferente. Desde luego, el analista debe tener en la mente las más bien considerables diferencias entre los tamaños muestrales en este caso.

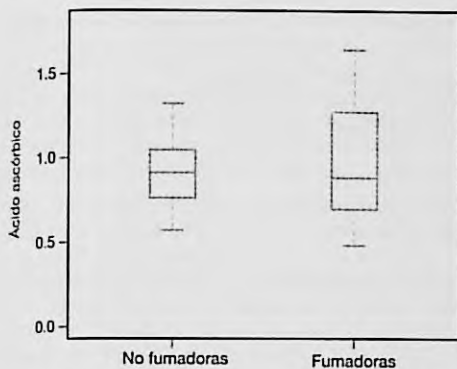


Figura 10.16: Dos gráficas de caja y bigote con los datos de ácido ascórbico para mujeres fumadoras y no fumadoras.



Figura 10.17: Dos gráficas de caja y bigote para los datos de los tallos.

Considere el ejercicio 9.40 de la sección 9.9. En la figura 10.17 se presenta la gráfica múltiple de caja y bigote para los datos de 10 tallos, de los cuales sólo la mitad recibió el tratamiento con nitrógeno. Tal gráfica revela una variabilidad menor para el grupo que no recibió nitrógeno. Además, la falta de traslape de las cajas sugiere una diferencia significativa entre los pesos medios de los tallos para los dos grupos. Parecería que la presencia de nitrógeno aumenta el peso de los tallos y quizás aumente la variabilidad en los pesos.

No existen reglas generales relacionadas con el momento cuando dos gráficas de caja y bigote brindan evidencia de diferencias significativas entre las medias. Sin embargo, una pauta aproximada es que si la línea del percentil 25 para una muestra excede a la línea de la mediana de la otra muestra, hay evidencia sólida de una diferencia entre las medias.

Se hará más énfasis en los métodos gráficos en un estudio de caso de la vida real que se presenta más adelante en este capítulo.

Impresiones por computadora con comentarios para pruebas t con dos muestras

Considere nuevamente el ejercicio 9.40 de la página 294, donde se reunieron datos de tallos que recibieron y no recibieron nitrógeno. Pruebe

$$H_0: \mu_{\text{NT}} = \mu_{\text{NO}},$$

$$H_1: \mu_{\text{NT}} > \mu_{\text{NO}},$$

donde las medias de la población indican los pesos medios. La figura 10.18 es una impresión por computadora con comentarios generados con el programa SAS. Observe que se presentan la desviación estándar y el error estándar muestrales para ambas muestras. También se incluye el estadístico t bajo la suposición de varianzas iguales y varianzas diferentes. En la gráfica de caja y bigote que se observa en la figura 10.17 en realidad parece que se transgrede la suposición de igualdad de varianzas. Un valor P de 0.0229 sugiere una conclusión de medias diferentes. Esto coincide con la información de diagnóstico que se presenta en la figura 10.18. A propósito, observe que t y t' son iguales en este caso, ya que $n_1 = n_2$.

TTEST Procedure					
Variable Weight					
Mineral	N	Mean	Std Dev	Std Err	
No nitrogen	10	0.3990	0.0728	0.0230	
Nitrogen	10	0.5650	0.1867	0.0591	
Variances		DF	t Value	Pr > t	
Equal	18	2.62	0.0174		
Unequal	11.7	2.62	0.0229		
Test the Equality of Variances					
Variable	Num DF	Den DF	F Value	Pr > F	
Weight	9	9	6.58	0.0098	

Figura 10.18: Impresión del SAS para la prueba t de dos muestras.

Ejercicios

10.19 En un informe de investigación, Richard H. Weindruch, de la Escuela de Medicina de la UCLA, afirma que los ratones con una vida promedio de 32 meses vivirán hasta alrededor de 40 meses si 40% de las calorías en su dieta se reemplazan con vitaminas y proteínas. ¿Hay alguna razón para creer que $\mu < 40$, si 64 ratones que son sometidos a esa dieta tienen una vida promedio de 38 meses, con una desviación estándar de 5.8 meses? Utilice un valor P en su conclusión.

10.20 Una muestra aleatoria de 64 bolsas de palomitas con queso chedar pesan, en promedio, 5.23 onzas, con una desviación estándar de 0.24 onzas. Pruebe la hipótesis de que $\mu = 5.5$ onzas contra la hipótesis alternativa de que $\mu < 5.5$ onzas, al nivel de significancia de 0.05.

10.21 Una empresa de material eléctrico fabrica bombillas que tienen una duración que se distribuye de forma aproximadamente normal con una media de 800 horas y una desviación estándar de 40 horas. Pruebe la hipótesis de que $\mu = 800$ horas contra la alternativa de que $\mu \neq 800$ horas, si una muestra aleatoria de 30 bombillas tiene una duración promedio de 788 horas. Utilice un valor P en su respuesta.

10.22 En la revista *Hypertension* de la American Heart Association, investigadores reportan que los individuos que practican la meditación trascendental (MT) bajan su presión sanguínea de forma significativa. Si una muestra aleatoria de 225 hombres que practican la MT meditan 8.5 horas a la semana, con una desviación estándar de 2.25 horas, ¿esto sugiere que, en promedio, los hombres que utilizan la MT meditan más de 8 horas por semana? Cite un valor P en su conclusión.

10.23 Pruebe la hipótesis de que el contenido promedio de los envases de un lubricante específico es de 10 litros, si los contenidos de una muestra aleatoria de 10 envases son: 10.2, 9.7, 10.1, 10.3, 10.1, 9.8, 9.9, 10.4, 10.3 y 9.8 litros. Utilice un nivel de significancia de 0.01

y suponga que la distribución del contenido es normal.

10.24 La estatura promedio de mujeres en el grupo de primer año de cierta universidad ha sido, históricamente, de 162.5 centímetros, con una desviación estándar de 6.9 centímetros. ¿Existe alguna razón para creer que ha habido un cambio en la estatura promedio, si una muestra aleatoria de 50 mujeres del grupo actual de primer año tiene una estatura promedio de 165.2 centímetros? Utilice un valor P en su conclusión. Suponga que la desviación estándar permanece constante.

10.25 Se afirma que los automóviles recorren en promedio más de 20,000 kilómetros por año. Para probar tal afirmación se pide a una muestra de 100 propietarios de automóviles seleccionada de manera aleatoria que lleven un registro de los kilómetros que recorren. ¿Estaría usted de acuerdo con esta afirmación, si la muestra aleatoria indicara un promedio de 23,500 kilómetros y una desviación estándar de 3900 kilómetros? Utilice un valor P en su conclusión.

10.26 De acuerdo con un estudio sobre un régimen alimenticio, la ingesta elevada de sodio se relaciona con úlceras, cáncer estomacal y migrañas. El requerimiento humano de sal es de tan sólo 220 miligramos diarios, el cual se rebasa en la mayoría de las porciones individuales de cereales listos para comerse. Si una muestra aleatoria de 20 porciones similares de cierto cereal tiene un contenido medio de 244 miligramos de sodio y una desviación estándar de 24.5 miligramos, ¿esto sugiere, a un nivel de significancia de 0.05, que el contenido promedio de sodio para porciones individuales de ese cereal es mayor que 220 miligramos? Suponga que la distribución de contenidos de sodio es normal.

10.27 Un estudio de la Universidad de Colorado en Boulder revela que correr aumenta el porcentaje de la tasa metabólica basal (TMB) en mujeres ancianas. La TMB promedio de 30 ancianas corredoras fue 34.0%

más alta que la TMB promedio de 30 ancianas sedentarias, en tanto que las desviaciones estándar reportadas fueron de 10.5 y 10.2%, respectivamente. ¿Existe un aumento significativo en la TMB de las corredoras respecto a las sedentarias? Suponga que las poblaciones se distribuyen de forma aproximadamente normal con varianzas iguales. Utilice un valor P en sus conclusiones.

10.28 De acuerdo con *Chemical Engineering*, una propiedad importante de la fibra es su absorbencia de agua. Se encontró que el porcentaje promedio de absorción de 25 pedazos de fibra de algodón seleccionados al azar es 20, con una desviación estándar de 1.5. Una muestra aleatoria de 25 pedazos de acetato reveló un porcentaje promedio de 12 con una desviación estándar de 1.25. ¿Existe evidencia sólida de que el porcentaje promedio de absorción de la población es significativamente mayor para la fibra de algodón que para el acetato? Suponga que el porcentaje de absorbencia se distribuye de forma casi normal y que las varianzas de la población en el porcentaje de absorbencia para las dos fibras son iguales. Utilice un nivel de significancia de 0.05.

10.29 La experiencia indica que el tiempo que requieren los estudiantes de último año de preparatoria para contestar una prueba estandarizada es una variable aleatoria normal con una media de 35 minutos. Si a una muestra aleatoria de 20 estudiantes de último año de preparatoria le toma un promedio de 33.1 minutos contestar esa prueba con una desviación estándar de 4.3 minutos, pruebe la hipótesis de que, a un nivel de significancia de 0.05, $\mu = 35$ minutos, contra la alternativa de que $\mu < 35$ minutos.

10.30 Una muestra aleatoria de tamaño $n_1 = 25$, tomada de una población normal con una desviación estándar $\sigma_1 = 5.2$, tiene una media $\bar{x}_1 = 81$. Una segunda muestra aleatoria de tamaño $n_2 = 36$, que se toma de una población normal diferente con una desviación estándar $\sigma_2 = 3.4$, tiene una media $\bar{x}_2 = 76$. Pruebe la hipótesis de que $\mu_1 = \mu_2$ contra la alternativa $\mu_1 \neq \mu_2$. Cite un valor P en su conclusión.

10.31 Un fabricante afirma que la resistencia promedio a la tensión del hilo A excede a la resistencia a la tensión promedio del hilo B en al menos 12 kilogramos. Para probar esta afirmación se pusieron a prueba 50 pedazos de cada tipo de hilo en condiciones similares. El hilo tipo A tuvo una resistencia promedio a la tensión de 86.7 kilogramos con una desviación estándar de 6.28 kilogramos; mientras que el hilo tipo B tuvo una resistencia promedio a la tensión de 77.8 kilogramos con una desviación estándar de 5.61 kilogramos. Pruebe la afirmación del fabricante usando un nivel de significancia de 0.05.

10.32 El *Amstat News* (diciembre de 2004) lista los sueldos medios de profesores asociados de estadística en instituciones de investigación, en escuelas de huma-

nidades y en otras instituciones en Estados Unidos. Suponga que una muestra de 200 profesores asociados de instituciones de investigación tiene un sueldo promedio de \$70,750 anuales con una desviación estándar de \$6000. Suponga también que una muestra de 200 profesores asociados de otros tipos de instituciones tienen un sueldo promedio de \$65,200 con una desviación estándar de \$5000. Pruebe la hipótesis de que el sueldo medio de profesores asociados de instituciones de investigación es \$2000 más alto que el de los profesores de otras instituciones. Utilice un nivel de significancia de 0.01.

10.33 Se llevó a cabo un estudio para saber si el aumento en la concentración de sustrato tiene un efecto apreciable sobre la velocidad de una reacción química. Con una concentración de sustrato de 1.5 moles por litro, la reacción se realizó 15 veces, con una velocidad promedio de 7.5 micromoles por 30 minutos y una desviación estándar de 1.5. Con una concentración de sustrato de 2.0 moles por litro, se realizaron 12 reacciones que produjeron una velocidad promedio de 8.8 micromoles por 30 minutos y una desviación estándar muestral de 1.2. ¿Hay alguna razón para creer que este incremento en la concentración de sustrato ocasiona un aumento en la velocidad media de la reacción de más de 0.5 micromoles por 30 minutos? Utilice un nivel de significancia de 0.01 y suponga que las poblaciones se distribuyen de forma aproximadamente normal con varianzas iguales.

10.34 Se realizó un estudio para determinar si los temas de un curso de física se comprenden mejor cuando éste incluye prácticas de laboratorio. Se seleccionaron estudiantes al azar para que participaran en un curso de tres semestres con una hora de clase sin prácticas de laboratorio o en un curso de cuatro semestres con una hora de clase con prácticas de laboratorio. En la sección con prácticas de laboratorio 11 estudiantes obtuvieron una calificación promedio de 85 con una desviación estándar de 4.7; mientras que en la sección sin prácticas de laboratorio 17 estudiantes obtuvieron una calificación promedio de 79 con una desviación estándar de 6.1. ¿Diría usted que el curso que incluyó prácticas de laboratorio aumentó la calificación promedio hasta en 8 puntos? Utilice un valor P en su conclusión y suponga que las poblaciones se distribuyen de forma aproximadamente normal con varianzas iguales.

10.35 Para indagar si un nuevo suero frena el desarrollo de la leucemia se seleccionan 9 ratones, todos en una etapa avanzada de la enfermedad. Cinco ratones reciben el tratamiento y cuatro no. Los tiempos de supervivencia, en años, a partir del momento en que comienza el experimento son los siguientes:

Con tratamiento	2.1	5.3	1.4	4.6	0.9
Sin tratamiento	1.9	0.5	2.8	3.1	

A un nivel de significancia de 0.05, ¿se puede decir que el suero es eficaz? Suponga que las dos poblaciones se distribuyen de forma normal con varianzas iguales.

10.36 Los ingenieros de una armadora de automóviles de gran tamaño están tratando de decidir si comprarán neumáticos de la marca *A* o de la marca *B* para sus modelos nuevos. Con el fin de ayudarlos a tomar una decisión se realiza un experimento en el que se usan 12 neumáticos de cada marca. Los neumáticos se utilizan hasta que se desgastan. Los resultados son los siguientes:

Marca *A*:

$$\bar{x}_1 = 37,900 \text{ kilómetros,}$$

$$s_1 = 5100 \text{ kilómetros.}$$

Marca *B*:

$$\bar{x}_2 = 39,800 \text{ kilómetros,}$$

$$s_2 = 5900 \text{ kilómetros.}$$

Pruebe la hipótesis de que no hay diferencia en el desgaste promedio de las 2 marcas de neumáticos. Suponga que las poblaciones se distribuyen de forma aproximadamente normal con varianzas iguales. Use un valor *P*.

10.37 En el ejercicio 9.42 de la página 295 pruebe la hipótesis de que el ahorro de combustible de los camiones compactos Volkswagen, en promedio, excede al de los camiones compactos Toyota equipados de forma similar, que utilizan 4 kilómetros por litro. Utilice un nivel de significancia de 0.10.

10.38 Un investigador de la UCLA afirma que el promedio de vida de los ratones se puede prolongar hasta por 8 meses cuando se reducen las calorías en su dieta aproximadamente 40% desde el momento en que se destetan. Las dietas restringidas se enriquecen a niveles normales con vitaminas y proteínas. Suponga que a una muestra aleatoria de 10 ratones que tienen una vida promedio de 32.1 meses con una desviación estándar de 3.2 meses se les alimenta con una dieta normal, mientras que a una muestra aleatoria de 15 ratones que tienen un promedio de vida de 37.6 meses con una desviación estándar de 2.8 meses se les alimenta con la dieta restringida. A un nivel de significancia de 0.05 pruebe la hipótesis de que el promedio de vida de los ratones con esta dieta restringida aumenta 8 meses, contra la alternativa de que el aumento es menor de 8 meses. Suponga que las distribuciones de la esperanza de vida con las dietas regular y restringida son aproximadamente normales con varianzas iguales.

10.39 Los siguientes datos representan los tiempos de duración de películas producidas por 2 empresas cinematográficas:

Empresa	Tiempo (minutos)					
1	102	86	98	109	92	
2	81	165	97	134	92	114

Pruebe la hipótesis de que la duración promedio de las películas producidas por la empresa 2 excede al tiempo promedio de duración de las que produce la empresa 1 en 10 minutos, contra la alternativa unilateral de que la diferencia es de menos de 10 minutos. Utilice un nivel de significancia de 0.1 y suponga que las distribuciones de la duración son aproximadamente normales con varianzas iguales.

10.40 En un estudio realizado en Virginia Tech se compararon los niveles de ácido ascórbico en plasma en mujeres embarazadas fumadoras con los de mujeres no fumadoras. Para el estudio se seleccionaron 32 mujeres que estuvieran en los últimos 3 meses de embarazo, que no tuvieran padecimientos importantes y que sus edades fluctuaran entre los 15 y los 32 años. Antes de tomar muestras de 20 ml de sangre se pidió a las participantes que fueran en ayunas, que no tomaran sus suplementos vitamínicos y que evitaran alimentos con alto contenido de ácido ascórbico. A partir de las muestras de sangre se determinaron los siguientes valores de ácido ascórbico en el plasma de cada mujer, en miligramos por 100 mililitros:

Valores de ácido ascórbico en plasma

No fumadoras	Fumadoras	
0.97	1.16	0.48
0.72	0.86	0.71
1.00	0.85	0.98
0.81	0.58	0.68
0.62	0.57	1.18
1.32	0.64	1.36
1.24	0.98	0.78
0.99	1.09	1.64
0.90	0.92	
0.74	0.78	
0.88	1.24	
0.94	1.18	

¿Existe suficiente evidencia para concluir que hay una diferencia entre los niveles de ácido ascórbico en plasma de mujeres fumadoras y no fumadoras? Suponga que los dos conjuntos de datos provienen de poblaciones normales con varianzas diferentes. Utilice un valor *P*.

10.41 El Departamento de Zoología de Virginia Tech llevó a cabo un estudio para determinar si existe una diferencia significativa en la densidad de organismos en dos estaciones diferentes ubicadas en Cedar Run, una corriente secundaria que se localiza en la cuenca del río Roanoke. El drenaje de una planta de tratamiento de aguas negras y el sobreflujo del estanque de sedimentación de la Federal Mogul Corporation entran al flujo cerca del nacimiento del río. Los siguientes datos proporcionan las medidas de densidad, en número de organismos por metro cuadrado, en las dos estaciones colectoras:

Número de organismos por metro cuadrado

Estación 1		Estación 2	
5030	4980	2800	2810
13,700	11,910	4670	1330
10,730	8130	6890	3320
11,400	26,850	7720	1230
860	17,660	7030	2130
2200	22,800	7330	2190
4250	1130		
15,040	1690		

A un nivel de significancia de 0.05, ¿podemos concluir que las densidades promedio en las dos estaciones son iguales? Suponga que las observaciones provienen de poblaciones normales con varianzas diferentes.

10.42 Cinco muestras de una sustancia ferrosa se usaron para determinar si existe una diferencia entre un análisis químico de laboratorio y un análisis de fluorescencia de rayos X del contenido de hierro. Cada muestra se dividió en dos submuestras y se aplicaron los dos tipos de análisis. A continuación se presentan los datos codificados que muestran los análisis de contenido de hierro:

Análisis	Muestra				
	1	2	3	4	5
Rayos X	2.0	2.0	2.3	2.1	2.4
Químico	2.2	1.9	2.5	2.3	2.4

Suponga que las poblaciones son normales y pruebe, al nivel de significancia de 0.05, si los dos métodos de análisis dan, en promedio, el mismo resultado.

10.43 De acuerdo con informes publicados, el ejercicio en condiciones de fatiga altera los mecanismos que determinan el desempeño. Se realizó un experimento con 15 estudiantes universitarios hombres, entrenados para realizar un movimiento horizontal continuo del brazo, de derecha a izquierda, desde un microinterruptor hasta una barrera, golpeando sobre la barrera en coincidencia con la llegada de una manecilla del reloj a la posición de las 6 en punto. Se registró el valor absoluto de la diferencia entre el tiempo, en milisegundos, que toma golpear sobre la barrera y el tiempo para que la manecilla alcance la posición de las 6 en punto (500 mseg). Cada participante ejecutó la tarea cinco veces en condiciones sin fatiga y con fatiga, y se registraron las siguientes sumas de las diferencias absolutas para las cinco ejecuciones:

Sujeto	Diferencias absolutas de tiempo	
	Sin fatiga	Con fatiga
1	158	91
2	92	59
3	65	215
4	98	226
5	33	223
6	89	91
7	148	92
8	58	177
9	142	134
10	117	116

11	74	153
12	66	219
13	109	143
14	57	164
15	85	100

Un aumento en la diferencia media absoluta de tiempo cuando la tarea se ejecuta en condiciones de fatiga apoyaría la afirmación de que el ejercicio, en condiciones de fatiga, altera el mecanismo que determina el desempeño. Suponga que las poblaciones se distribuyen normalmente y pruebe tal afirmación.

10.44 En un estudio realizado por el Departamento de Nutrición Humana y Alimentos del Virginia Tech se registraron los siguientes datos sobre los residuos de ácido sórbico en jamón, en partes por millón, inmediatamente después de sumergirlo en una solución de sorbato y después de 60 días de almacenamiento:

Residuos de ácido sórbico en jamón

Rebanada	Antes del almacenamiento	Después del almacenamiento
1	224	116
2	270	96
3	400	239
4	444	329
5	590	437
6	660	597
7	1400	689
8	680	576

Si se supone que las poblaciones se distribuyen normalmente, ¿hay suficiente evidencia, a un nivel de significancia de 0.05, para decir que la duración del almacenamiento influye en las concentraciones residuales de ácido sórbico?

10.45 El administrador de una empresa de taxis está tratando de decidir si el uso de neumáticos radiales en lugar de neumáticos regulares cinturados mejora el rendimiento de combustible. Se equipan 12 autos con neumáticos radiales y se conducen en un recorrido de prueba preestablecido. Sin cambiar a los conductores, los mismos autos se equipan con neumáticos regulares cinturados y se conducen nuevamente en el recorrido de prueba. Se registraron los siguientes datos sobre el consumo de gasolina, en kilómetros por litro:

Kilómetros por litro

Automóvil	Llantas radiales	Llantas cinturadas
1	4.2	4.1
2	4.7	4.9
3	6.6	6.2
4	7.0	6.9
5	6.7	6.8
6	4.5	4.4
7	5.7	5.7
8	6.0	5.8
9	7.4	6.9
10	4.9	4.7
11	6.1	6.0
12	5.2	4.9

¿Podemos concluir que los autos equipados con neumáticos radiales ahorran más combustible que aquellos equipados con neumáticos cinturados? Suponga que las poblaciones se distribuyen normalmente. Utilice un valor P en su conclusión.

10.46 En el ejercicio de repaso 9.91 de la página 313 utilice la distribución t para probar la hipótesis de que la dieta reduce el peso de un individuo en 4.5 kilogramos, en promedio, contra la hipótesis alternativa de que la diferencia media en peso es menor que 4.5 kilogramos. Utilice un valor P .

10.47 ¿Qué tan grande debería ser la muestra del ejercicio 10.20 para que la potencia de la prueba sea de 0.90, cuando la media verdadera es 5.20? Suponga que $\sigma = 0.24$.

10.48 Si la distribución del tiempo de vida en el ejercicio 10.19 es aproximadamente normal, ¿qué tan grande debería ser una muestra para que la probabilidad de cometer un error tipo II sea 0.1 cuando la media verdadera es 35.9 meses? Suponga que $\sigma = 5.8$ meses.

10.49 ¿Qué tan grande debería ser la muestra del ejercicio 10.24 para que la potencia de la prueba sea de 0.95 cuando la estatura promedio verdadera difiere de 162.5 en 3.1 centímetros? Utilice $\alpha = 0.02$.

10.50 ¿Qué tan grandes deberían ser las muestras del ejercicio 10.31 para que la potencia de la prueba sea de 0.95, cuando la diferencia verdadera entre los tipos de hilo A y B es 8 kilogramos?

10.51 ¿Qué tan grande debería ser la muestra del ejercicio 10.22 para que la potencia de la prueba sea de 0.8 cuando el tiempo promedio verdadero dedicado a la meditación excede al valor hipotético en 1.2 σ ? Utilice $\alpha = 0.05$.

10.52 Se considera una prueba t a un nivel $\alpha = 0.05$ para probar

$$H_0: \mu = 14,$$

$$H_1: \mu \neq 14.$$

¿Qué tamaño de muestra se necesita para que la probabilidad de no rechazar de manera errónea H_0 sea 0.1 cuando la media de la población verdadera difiere de 14 en 0.5? A partir de una muestra preliminar estimamos que σ es 1.25.

10.53 En el Departamento de Medicina Veterinaria del Virginia Tech se llevó a cabo un estudio para determinar si la "resistencia" de una herida de incisión quirúrgica es afectada por la temperatura del bisturí. En el experimento se utilizaron 8 perros. Se hicieron incisiones "calientes" y "frías" en el abdomen de cada

perro y se midió la resistencia. A continuación se presentan los datos resultantes.

Perro	Bisturí	Resistencia
1	Caliente	5120
1	Frío	8200
2	Caliente	10,000
2	Frío	8600
3	Caliente	10,000
3	Frío	9200
4	Caliente	10,000
4	Frío	6200
5	Caliente	10,000
5	Frío	10,000
6	Caliente	7900
6	Frío	5200
7	Caliente	510
7	Frío	885
8	Caliente	1020
8	Frío	460

- Escriba una hipótesis adecuada para determinar si la resistencia de las incisiones realizadas con bisturí caliente difiere en forma significativa de la resistencia de las realizadas con bisturí frío.
- Pruebe la hipótesis utilizando una prueba t pareada. Utilice un valor P en su conclusión.

10.54 Se utilizaron 9 sujetos en un experimento para determinar si la exposición a monóxido de carbono tiene un impacto sobre la capacidad respiratoria. Los datos fueron recolectados por el personal del Departamento de Salud y Educación Física del Virginia Tech y analizados en el Centro de Consulta Estadística en Hokie Land. Los sujetos fueron expuestos a cámaras de respiración, una de las cuales contenía una alta concentración de CO. Se realizaron varias mediciones de frecuencia respiratoria a cada sujeto en cada cámara. Los sujetos fueron expuestos a las cámaras de respiración en una secuencia aleatoria. Los siguientes datos representan la frecuencia respiratoria en número de respiraciones por minuto. Realice una prueba unilateral de la hipótesis de que la frecuencia respiratoria media es igual en los dos ambientes. Utilice $\alpha = 0.05$. Suponga que la frecuencia respiratoria es aproximadamente normal.

Sujeto	Con CO	Sin CO
1	30	30
2	45	40
3	26	25
4	25	23
5	34	30
6	51	49
7	46	41
8	32	35
9	30	28

10.8 Una muestra: prueba sobre una sola proporción

Las pruebas de hipótesis que se relacionan con proporciones se requieren en muchas áreas. A los políticos les interesa conocer la fracción de votantes que los favorecerá en la siguiente elección. Todas las empresas manufactureras se preocupan por la proporción de artículos defectuosos cuando se realiza un embarque. Los jugadores dependen del conocimiento de la proporción de resultados que consideran favorables.

Consideraremos el problema de probar la hipótesis de que la proporción de éxitos en un experimento binomial es igual a algún valor específico. Es decir, probaremos la hipótesis nula H_0 de que $p = p_0$, donde p es el parámetro de la distribución binomial. La hipótesis alternativa puede ser una de las alternativas unilaterales o bilaterales usuales:

$$p < p_0, \quad p > p_0, \quad \text{o} \quad p \neq p_0.$$

La variable aleatoria adecuada sobre la que basamos nuestro criterio de decisión es la variable aleatoria binomial X ; aunque también podríamos usar el estadístico $\hat{p} = X/n$. Los valores de X que están lejos de la media $\mu = np_0$ conducirán al rechazo de la hipótesis nula. Como X es una variable binomial discreta, es poco probable que se pueda establecer una región crítica cuyo tamaño sea *exactamente* igual a un valor preestablecido de α . Por esta razón es preferible, al trabajar con muestras pequeñas, basar nuestras decisiones en valores P . Para probar la hipótesis

$$H_0: p = p_0,$$

$$H_1: p < p_0,$$

utilizamos la distribución binomial para calcular el valor P

$$P = P(X \leq x \text{ cuando } p = p_0).$$

El valor x es el número de éxitos en nuestra muestra de tamaño n . Si este valor P es menor o igual que α , nuestra prueba es significativa al nivel α y rechazamos H_0 a favor de H_1 . De manera similar, para probar la hipótesis

$$H_0: p = p_0,$$

$$H_1: p > p_0,$$

al nivel de significancia α , calculamos

$$P = P(X \geq x \text{ cuando } p = p_0)$$

y rechazamos H_0 a favor de H_1 si este valor P es menor o igual que α . Finalmente, para probar la hipótesis

$$H_0: p = p_0,$$

$$H_1: p \neq p_0,$$

a un nivel de significancia α , calculamos

$$P = 2P(X \leq x \text{ cuando } p = p_0) \quad \text{si } x < np_0$$

o

$$P = 2P(X \geq x \text{ cuando } p = p_0) \quad \text{si } x > np_0$$

y rechazamos H_0 a favor de H_1 si el valor P calculado es menor o igual que α .

Los pasos para probar una hipótesis nula acerca de una proporción contra varias alternativas usando las probabilidades binomiales de la tabla A.1 son los siguientes:

Prueba de una proporción (muestras pequeñas)	<ol style="list-style-type: none"> 1. $H_0: p = p_0$. 2. Una de las alternativas $H_1: p < p_0, p > p_0$ o $p \neq p_0$. 3. Elegir un nivel de significancia igual a α. 4. Estadístico de prueba: variable binomial X con $p = p_0$. 5. Cálculos: obtener x, el número de éxitos, y calcular el valor P adecuado. 6. Decisión: sacar las conclusiones apropiadas con base en el valor P.
--	--

Ejemplo 10.9: Un constructor afirma que en 70% de las viviendas que se construyen actualmente en la ciudad de Richmond, Virginia, se instalan bombas de calor. ¿Estaría de acuerdo con esta afirmación si una encuesta aleatoria de viviendas nuevas en esta ciudad revelara que 8 de 15 tienen instaladas bombas de calor? Utilice un nivel de significancia de 0.10.

- Solución:**
1. $H_0: p = 0.7$.
 2. $H_1: p \neq 0.7$.
 3. $\alpha = 0.10$.
 4. Estadístico de prueba: Variable binomial X con $p = 0.7$ y $n = 15$.
 5. Cálculos: $x = 8$ y $np_0 = (15)(0.7) = 10.5$. Por lo tanto, de la tabla A.1, el valor P calculado es

$$P = 2P(X \leq 8 \text{ cuando } p = 0.7) = 2 \sum_{x=0}^8 b(x; 15, 0.7) = 0.2622 > 0.10.$$

6. Decisión: No rechazar H_0 . Concluir que no hay razón suficiente para dudar de la afirmación del constructor. ▮

En la sección 5.2 aprendimos que cuando n es pequeña las probabilidades binomiales se pueden obtener de la fórmula binomial real o de la tabla A.1. Para n grande se requieren procedimientos de aproximación. Cuando el valor hipotético p_0 está muy cerca de 0 o de 1 se puede utilizar la distribución de Poisson con parámetro $\mu = np_0$. Sin embargo, para n grande por lo general se prefiere la aproximación de la curva normal, con los parámetros $\mu = np_0$ y $\sigma^2 = np_0q_0$, la cual es muy precisa, siempre y cuando p_0 no esté demasiado cerca de 0 o de 1. Si utilizamos la aproximación normal, el valor z para probar $p = p_0$ es dado por

$$z = \frac{x - np_0}{\sqrt{np_0q_0}} = \frac{\hat{p} - p_0}{\sqrt{p_0q_0/n}},$$

que es un valor de la variable normal estándar Z . Por consiguiente, para una prueba de dos colas al nivel de significancia α , la región crítica es $z < -z_{\alpha/2}$ o $z > z_{\alpha/2}$. Para la alternativa unilateral $p < p_0$, la región crítica es $z < -z_\alpha$, y para la alternativa $p > p_0$, la región crítica es $z > z_\alpha$.

Ejemplo 10.10: Se considera que un medicamento que se prescribe comúnmente para aliviar la tensión nerviosa tiene una eficacia de tan sólo 60%. Los resultados experimentales de un nuevo fármaco administrado a una muestra aleatoria de 100 adultos que padecían tensión nerviosa revelaron que 70 de ellos sintieron alivio. ¿Esta evidencia es suficiente para concluir que el nuevo medicamento es mejor que el que se prescribe comúnmente? Utilice un nivel de significancia de 0.05.

Solución: 1. $H_0: p = 0.6$.

2. $H_1: p > 0.6$.

3. $\alpha = 0.05$.

4. Región crítica: $z > 1.645$.

5. Cálculos: $x = 70$, $n = 100$, $\hat{p} = 70/100 = 0.7$ y

$$z = \frac{0.7 - 0.6}{\sqrt{(0.6)(0.4)/100}} = 2.04, P = P(Z > 2.04) < 0.0207.$$

6. Decisión: Rechazar H_0 y concluir que el nuevo fármaco es mejor.

10.9 Dos muestras: pruebas sobre dos proporciones

A menudo surgen situaciones en las que se desea probar la hipótesis de que dos proporciones son iguales. Por ejemplo, podemos tratar de mostrar evidencia de que la proporción de médicos que son pediatras en un estado es igual a la proporción de pediatras en otro estado. Quizás un individuo decida dejar de fumar sólo si se convence de que la proporción de fumadores con cáncer pulmonar excede a la proporción de no fumadores con ese tipo de cáncer.

En general, deseamos probar la hipótesis nula de que dos proporciones, o parámetros binomiales, son iguales. Es decir, probamos $p_1 = p_2$ contra una de las alternativas $p_1 < p_2$, $p_1 > p_2$, o $p_1 \neq p_2$. Desde luego, esto es equivalente a probar la hipótesis nula de que $p_1 - p_2 = 0$ contra una de las alternativas $p_1 - p_2 < 0$, $p_1 - p_2 > 0$ o $p_1 - p_2 \neq 0$. El estadístico sobre el que basamos nuestra decisión es la variable aleatoria $\hat{p}_1 - \hat{p}_2$. Se seleccionan al azar muestras independientes de tamaños n_1 y n_2 de dos poblaciones binomiales y se calcula la proporción de éxitos \hat{p}_1 y \hat{p}_2 para las dos muestras.

En la construcción de intervalos de confianza para p_1 y p_2 observamos, para n_1 y n_2 suficientemente grandes, que el estimador puntual \hat{p}_1 menos \hat{p}_2 estaba distribuido de forma casi normal con media

$$\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$$

y varianza

$$\sigma_{\hat{p}_1 - \hat{p}_2}^2 = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}.$$

Por lo tanto, es posible establecer la(s) región(es) crítica(s) usando la variable normal estándar

$$Z = \frac{(\hat{P}_1 - \hat{P}_2) - (p_1 - p_2)}{\sqrt{p_1 q_1 / n_1 + p_2 q_2 / n_2}}$$

Cuando H_0 es verdadera, podemos sustituir $p_1 = p_2 = p$ y $q_1 = q_2 = q$ (donde p y q son los valores comunes) en la fórmula anterior para Z y obtener la forma

$$Z = \frac{\hat{P}_1 - \hat{P}_2}{\sqrt{pq(1/n_1 + 1/n_2)}}$$

Sin embargo, para calcular un valor de Z debemos estimar los parámetros p y q que aparecen en el radical. Al agrupar los datos de ambas muestras el **estimado agrupado de la proporción** p es

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2},$$

donde x_1 y x_2 son el número de éxitos en cada una de las dos muestras. Al sustituir \hat{p} por p y $\hat{q} = 1 - \hat{p}$ por q , el valor z para probar $p_1 = p_2$ se determina a partir de la fórmula

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}(1/n_1 + 1/n_2)}}$$

Las regiones críticas para las hipótesis alternativas adecuadas se establecen como antes, utilizando puntos críticos de la curva normal estándar. En consecuencia, para la alternativa $p_1 \neq p_2$, al nivel de significancia α , la región crítica es $z < -z_{\alpha/2}$ o $z > z_{\alpha/2}$. Para una prueba donde la alternativa es $p_1 < p_2$, la región crítica será $z < -z_{\alpha}$; y cuando la alternativa es $p_1 > p_2$, la región crítica será $z > z_{\alpha}$.

Ejemplo 10.11: Se organizará una votación entre los residentes de una ciudad y el condado circundante para determinar si se aprueba una propuesta para la construcción de una planta química. Como el lugar en el que se propone construirla está dentro de los límites de la ciudad, muchos votantes del condado consideran que la propuesta será aprobada debido a la gran proporción de votantes que está a favor de que se construya. Se realiza una encuesta para determinar si hay una diferencia significativa en la proporción de votantes de la ciudad y los votantes del condado que favorecen la propuesta. Si 120 de 200 votantes de la ciudad favorecen la propuesta y 240 de 500 residentes del condado también lo hacen, ¿estaría usted de acuerdo en que la proporción de votantes de la ciudad que favorecen la propuesta es mayor que la proporción de votantes del condado? Utilice un nivel de significancia de $\alpha = 0.05$.

Solución: Sean p_1 y p_2 las proporciones verdaderas de votantes en la ciudad y el condado, respectivamente, que favorecen la propuesta.

1. $H_0: p_1 = p_2$.
2. $H_1: p_1 > p_2$.
3. $\alpha = 0.05$
4. Región crítica: $z > 1.645$.
5. Cálculos:

$$\hat{p}_1 = \frac{x_1}{n_1} = \frac{120}{200} = 0.60, \quad \hat{p}_2 = \frac{x_2}{n_2} = \frac{240}{500} = 0.48, \quad y$$

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2} = \frac{120 + 240}{200 + 500} = 0.51.$$

Por lo tanto,

$$z = \frac{0.60 - 0.48}{\sqrt{(0.51)(0.49)(1/200 + 1/500)}} = 2.9,$$

$$P = P(Z > 2.9) = 0.0019.$$

6. Decisión: Rechazar H_0 y estar de acuerdo en que la proporción de votantes de la ciudad a favor de la propuesta es mayor que la proporción de votantes del condado. ▮

Ejercicios

10.55 Un experto en mercadotecnia de una empresa fabricante de pasta considera que 40% de los amantes de la pasta prefieren la lasagna. Si 9 de 20 amantes de la pasta eligen la lasagna sobre otras pastas, ¿qué se puede concluir acerca de la afirmación del experto? Utilice un nivel de significancia de 0.05.

10.56 Suponga que, en el pasado, 40% de todos los adultos estaban a favor de la pena capital. ¿Existe alguna razón para creer que la proporción de adultos que está a favor de la pena capital ha aumentado si, en una muestra aleatoria de 15 adultos, 8 están a favor de la pena capital? Utilice un nivel de significancia de 0.05.

10.57 Se está considerando utilizar un nuevo aparato de radar para cierto sistema de misiles de defensa. El sistema se verifica experimentando con una aeronave en la que se simula una situación en la que alguien muere y otra en la que no ocurre ninguna muerte. Si en 300 ensayos ocurren 250 muertes, al nivel de significancia de 0.04, acepte o rechace la afirmación de que la probabilidad de una muerte con el nuevo sistema no excede a la probabilidad de 0.8 del sistema que se utiliza actualmente.

10.58 Se cree que al menos 60% de los residentes de cierta área están a favor de una demanda de anexión de una ciudad vecina. ¿Qué conclusión extraerá si sólo 110 de una muestra de 200 votantes están a favor de la demanda? Utilice un nivel de significancia de 0.05.

10.59 Una empresa petrolera afirma que en una quinta parte de las viviendas de cierta ciudad la gente utiliza petróleo como combustible para calentarlas. ¿Existen razones para creer que en menos de una quinta parte de las viviendas la gente utiliza este combustible para calentarlas si, en una muestra aleatoria

de 1000 viviendas de esa ciudad, se encuentra que 136 utilizan petróleo como combustible? Utilice un valor P en su conclusión.

10.60 En cierta universidad se estima que a lo sumo 25% de los estudiantes van en bicicleta a la escuela. ¿Parece que ésta es una estimación válida si, en una muestra aleatoria de 90 estudiantes universitarios, se encuentra que 28 van en bicicleta a la escuela? Utilice un nivel de significancia de 0.05.

10.61 En un invierno con epidemia de influenza los investigadores de una conocida empresa farmacéutica encuestaron a los padres de 2000 bebés para determinar si el nuevo medicamento de la empresa era eficaz después de dos días. De 120 bebés que tenían influenza y que recibieron el medicamento, 29 se curaron en dos días o menos. De 280 bebés que tenían influenza pero no recibieron el fármaco, 56 se curaron en dos días o menos. ¿Hay alguna indicación significativa que apoye la afirmación de la empresa sobre la eficacia del medicamento?

10.62 En un experimento de laboratorio controlado, científicos de la Universidad de Minnesota descubrieron que 25% de cierta cepa de ratas sujetas a una dieta con 20% de grano de café y luego forzadas a consumir un poderoso químico causante de cáncer desarrollaron tumores cancerosos. Si el experimento se repite, y 16 de 48 ratas desarrollan tumores, ¿existen razones para creer que la proporción de ratas que desarrollan tumores cuando se someten a esta dieta se incrementa? Utilice un nivel de significancia de 0.05.

10.63 En un estudio que se realizó para estimar la proporción de residentes de cierta ciudad y sus suburbios que están a favor de que se construya una planta

de energía nuclear se encontró que 63 de 100 residentes urbanos están a favor de la construcción, mientras que sólo 59 de 125 residentes suburbanos la apoyan. ¿Hay una diferencia significativa entre la proporción de residentes urbanos y suburbanos que están a favor de que se construya la planta nuclear? Utilice un valor P .

10.64 En un estudio sobre la fertilidad de mujeres casadas, realizado por Martin O'Connell y Carolyn C. Rogers para la Oficina del Censo en 1979, se seleccionaron al azar dos grupos de mujeres casadas de entre 25 y 29 años de edad y sin hijos, y a cada una se le preguntó si planeaba tener un hijo en algún momento. Se seleccionó un grupo de mujeres con menos de dos años de casadas y otro de mujeres con cinco años de casadas. Suponga que 240 de 300 mujeres con menos de dos años de casadas planean tener un hijo algún día, en comparación con 288 de las 400 mujeres con cinco años de casadas. ¿Podemos concluir que la proporción de mujeres con menos de dos años de casadas que planean tener hijos es significativamente mayor que la proporción de mujeres con cinco años de casadas que también planean tenerlos? Utilice un valor P .

10.65 Una comunidad urbana quiere demostrar que la incidencia de cáncer de mama es mayor en su localidad que en una área rural vecina. (Se encontró que los niveles de PCB son más altos en el suelo de la comunidad urbana). Si descubre que en la comunidad urbana 20 de 200 mujeres adultas tienen cáncer de mama y que en la comunidad rural 10 de 150 mujeres adultas lo tienen, ¿podría concluir, con un nivel de significancia de 0.05, que el cáncer de mama prevalece más en la comunidad urbana?

10.66 Proyecto de grupo: Para este proyecto el grupo se debe dividir en parejas. Suponga que se supone que al menos 25% de los estudiantes de su universidad hacen más de dos horas de ejercicio por semana. Reúna datos de una muestra aleatoria de 50 estudiantes y pregunte a cada uno si se ejercita durante al menos dos horas por semana; luego haga los cálculos necesarios para rechazar o no rechazar la suposición anterior. Demuestre todo el procedimiento y utilice un valor P en sus conclusiones.

10.10 Pruebas de una y dos muestras referentes a varianzas

En esta sección estudiaremos la prueba de hipótesis relacionada con varianzas o desviaciones estándar de la población. No son poco comunes las aplicaciones de pruebas de una y dos muestras sobre varianzas. Los ingenieros y los científicos constantemente se enfrentan a estudios donde se les pide demostrar que las mediciones que tienen que ver con productos o procesos cumplen con las especificaciones que fijan los consumidores. Las especificaciones a menudo se cumplen si la varianza del proceso es suficientemente pequeña. También existe interés por experimentos que comparan métodos o procesos donde la reproducibilidad o variabilidad inherentes se deben comparar de manera formal. Además, para determinar si no se cumple la suposición de varianzas iguales, con frecuencia se aplica una prueba que compara dos varianzas antes de llevar a cabo una prueba t sobre dos medias.

Empecemos por considerar el problema de probar la hipótesis nula H_0 de que la varianza de la población σ^2 es igual a un valor específico σ_0^2 contra una de las alternativas comunes $\sigma^2 < \sigma_0^2$, $\sigma^2 > \sigma_0^2$ o $\sigma^2 \neq \sigma_0^2$. El estadístico apropiado sobre el que basamos nuestra decisión es el estadístico chi cuadrada del teorema 8.4, el cual se utilizó en el capítulo 9 para construir un intervalo de confianza para σ^2 . Por lo tanto, si suponemos que la distribución de la población que se muestra es normal, el valor de chi cuadrada para probar $\sigma^2 = \sigma_0^2$ es dado por

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2},$$

donde n es el tamaño de la muestra, s^2 es la varianza muestral y σ_0^2 es el valor de σ^2 dado por la hipótesis nula. Si H_0 es verdadera, χ^2 es un valor de la distribución chi cuadrada con $v = n - 1$ grados de libertad. En consecuencia, para una prueba de dos colas a un

nivel de significancia α , la región crítica es $\chi^2 < \chi^2_{1-\alpha/2}$ o $\chi^2 > \chi^2_{\alpha/2}$. Para la alternativa unilateral $\sigma^2 < \sigma_0^2$, la región crítica es $\chi^2 < \chi^2_{1-\alpha}$; y para la alternativa unilateral $\sigma^2 > \sigma_0^2$, la región crítica es $\chi^2 > \chi^2_{\alpha}$.

Robustez de la prueba χ^2 para la suposición de normalidad

Tal vez el lector se habrá dado cuenta de que varias pruebas dependen, al menos en teoría, de la suposición de normalidad. En general muchos procedimientos en estadística aplicada tienen fundamentos teóricos que dependen de la distribución normal. Estos procedimientos varían en el grado en que dependen de la suposición de la normalidad. A un procedimiento que es razonablemente insensible a esta suposición se le denomina **procedimiento robusto**, es decir, robusto para la normalidad. La prueba χ^2 sobre una sola varianza no es robusta en absoluto para la normalidad, es decir, el éxito práctico del procedimiento depende de la normalidad. Como resultado, el valor P calculado podría ser notablemente diferente del valor P verdadero si la población de la que se toma la muestra no es normal. De hecho, resulta muy plausible que un valor P estadísticamente significativo no sea una verdadera señal de $H_1: \sigma \neq \sigma_0$, sino que un valor significativo sea el resultado de haber violado las suposiciones de normalidad. Por lo tanto, el analista debería utilizar esta prueba χ^2 específica con precaución.

Ejemplo 10.12: Un fabricante de baterías para automóvil afirma que la duración de sus baterías se distribuye de forma aproximadamente normal con una desviación estándar igual a 0.9 años. Si una muestra aleatoria de 10 de tales baterías tiene una desviación estándar de 1.2 años, ¿considera que $\sigma > 0.9$ años? Utilice un nivel de significancia de 0.05.

- Solución:**
1. $H_0: \sigma^2 = 0.81$.
 2. $H_1: \sigma^2 > 0.81$.
 3. $\alpha = 0.05$.
 4. Región crítica: En la figura 10.19 vemos que se rechaza la hipótesis nula cuando $\chi^2 > 16.919$, donde $\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$ con $v = 9$ grados de libertad.

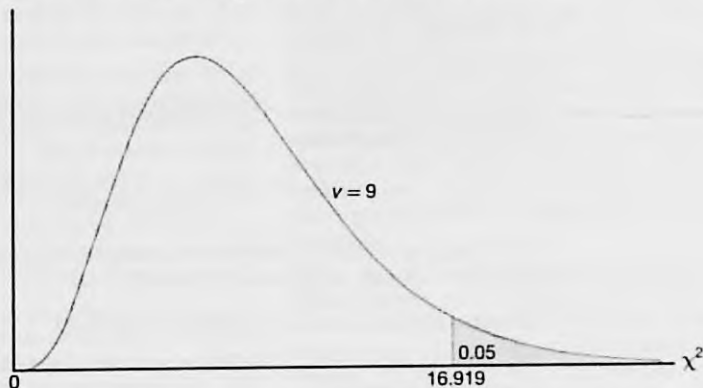


Figura 10.19: Región crítica para la hipótesis alternativa $\sigma > 0.9$.

5. Cálculos: $s^2 = 1.44$, $n = 10$ y

$$\chi^2 = \frac{(9)(1.44)}{0.81} = 16.0, \quad P \approx 0.07.$$

6. Decisión: El estadístico χ^2 no es significativo al nivel 0.05. Sin embargo, con base en el valor P de 0.07, hay evidencia de que $\sigma > 0.9$. ▮

Consideremos ahora el problema de probar la igualdad de las varianzas σ_1^2 y σ_2^2 de dos poblaciones. Esto es, probaremos la hipótesis nula H_0 de que $\sigma_1^2 = \sigma_2^2$ contra una de las alternativas usuales

$$\sigma_1^2 < \sigma_2^2, \quad \sigma_1^2 > \sigma_2^2, \quad \text{o} \quad \sigma_1^2 \neq \sigma_2^2.$$

Para muestras aleatorias independientes de tamaños n_1 y n_2 , respectivamente, de las dos poblaciones, el valor f para probar $\sigma_1^2 = \sigma_2^2$ es el cociente

$$f = \frac{s_1^2}{s_2^2},$$

donde s_1^2 y s_2^2 son las varianzas calculadas de las dos muestras. Si las dos poblaciones se distribuyen de forma aproximadamente normal y la hipótesis nula es verdadera, de acuerdo con el teorema 8.8 el cociente $f = s_1^2 / s_2^2$ es un valor de la distribución F con $\nu_1 = n_1 - 1$ y $\nu_2 = n_2 - 1$ grados de libertad. Por lo tanto, las regiones críticas de tamaño α que corresponden a las alternativas unilaterales $\sigma_1^2 < \sigma_2^2$ y $\sigma_1^2 > \sigma_2^2$ son, respectivamente, $f < f_{1-\alpha}(\nu_1, \nu_2)$ y $f > f_{\alpha}(\nu_1, \nu_2)$. Para la alternativa bilateral $\sigma_1^2 \neq \sigma_2^2$ la región crítica es $f < f_{1-\alpha/2}(\nu_1, \nu_2)$ o $f > f_{\alpha/2}(\nu_1, \nu_2)$.

Ejemplo 10.13: Al probar la diferencia en el desgaste abrasivo de los dos materiales del ejemplo 10.6 supusimos que las dos varianzas de la población desconocidas eran iguales. ¿Se justifica tal suposición? Utilice un nivel de significancia de 0.10.

Solución: Sean σ_1^2 y σ_2^2 las varianzas de la población para el desgaste abrasivo del material 1 y del material 2, respectivamente.

1. $H_0: \sigma_1^2 = \sigma_2^2$

2. $H_1: \sigma_1^2 \neq \sigma_2^2$

3. $\alpha = 0.10$.

4. Región crítica: En la figura 10.20 observamos que $f_{0.05}(11, 9) = 3.11$, y, usando el teorema 8.7, encontramos

$$f_{0.95}(11, 9) = \frac{1}{f_{0.05}(9, 11)} = 0.34.$$

Por lo tanto, se rechaza la hipótesis nula cuando $f < 0.34$ o $f > 3.11$, donde $f = s_1^2 / s_2^2$ con $\nu_1 = 11$ y $\nu_2 = 9$ grados de libertad.

5. Cálculos: $s_1^2 = 16$, $s_2^2 = 25$, por ende, $f = \frac{16}{25} = 0.64$.

6. Decisión: no rechazar H_0 . Concluir que no hay suficiente evidencia de que las varianzas sean diferentes. ▮

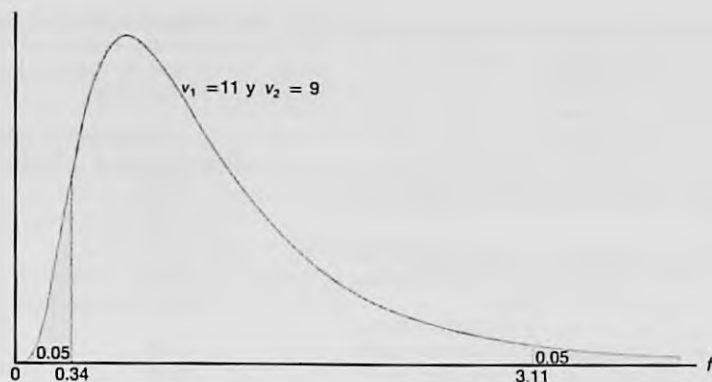


Figura 10.20: Región crítica para la hipótesis alternativa $\sigma_1^2 \neq \sigma_2^2$.

Prueba F para la prueba de varianzas con el SAS

La figura 10.18 de la página 356 presenta la impresión de una prueba t de dos muestras donde se comparan dos medias de los datos de los tallos en el ejercicio 9.40. La gráfica de caja y bigote que se observa en la figura 10.17 de la página 355 sugiere que las varianzas no son homogéneas y, por consiguiente, el estadístico t' y su valor P correspondiente son relevantes. Observe también que la impresión muestra el estadístico F para $H_0: \sigma_1 = \sigma_2$ con un valor P de 0.0098, que es evidencia adicional de que se debe esperar más variabilidad cuando se aplica el tratamiento con nitrógeno que cuando no se aplica.

Ejercicios

10.67 Se sabe que el contenido de los envases de un lubricante específico se distribuye normalmente con una varianza de 0.03 litros. Pruebe la hipótesis de que $\sigma^2 = 0.03$ contra la alternativa de que $\sigma^2 \neq 0.03$ para la muestra aleatoria de 10 envases del ejercicio 10.23 de la página 356. Use un valor P en sus conclusiones.

10.68 Por experiencia se sabe que el tiempo que se requiere para que los estudiantes de preparatoria de último año contesten una prueba estandarizada es una variable aleatoria normal con una desviación estándar de 6 minutos. Pruebe la hipótesis de que $\sigma = 6$ contra la alternativa de que $\sigma < 6$ si una muestra aleatoria de los tiempos para realizar la prueba de 20 estudiantes de preparatoria de último año tiene una desviación estándar $s = 4.51$. Utilice un nivel de significancia de 0.05.

10.69 Se deben supervisar las aflotoxinas ocasionadas por moho en cosechas de cacahuate en Virginia. Una muestra de 64 lotes de cacahuate revela niveles de 24.17 ppm, en promedio, con una varianza de 4.25 ppm. Pruebe la hipótesis de que $\sigma^2 = 4.2$ ppm contra la alternativa de que $\sigma^2 \neq 4.2$ ppm. Utilice un valor P en sus conclusiones.

10.70 Datos históricos indican que la cantidad de dinero que aportaron los residentes trabajadores de una ciudad grande para un escuadrón de rescate voluntario es una variable aleatoria normal con una desviación estándar de \$1.40. Se sugiere que las contribuciones al escuadrón de rescate sólo de los empleados del departamento de sanidad son mucho más variables. Si las contribuciones de una muestra aleatoria de 12 empleados del departamento de sanidad tienen una desviación estándar de \$1.75, ¿podemos concluir a un nivel de significancia de 0.01 que la desviación estándar de las contribuciones de todos los trabajadores de sanidad es mayor que la de todos los trabajadores que viven en dicha ciudad?

10.71 Se dice que una máquina despachadora de bebida gaseosa está fuera de control si la varianza de los contenidos excede a 1.15 decilitros. Si una muestra aleatoria de 25 bebidas de esta máquina tiene una varianza de 2.03 decilitros, ¿esto indica, a un nivel de significancia de 0.05, que la máquina está fuera de control? Suponga que los contenidos se distribuyen de forma aproximadamente normal.

10.72 Prueba de $\sigma^2 = \sigma_0^2$ para una muestra grande: Cuando $n \geq 30$ podemos probar la hipótesis nula de que $\sigma^2 = \sigma_0^2$ o $\sigma = \sigma_0$ calculando

$$z = \frac{s - \sigma_0}{\sigma_0 / \sqrt{2n}},$$

que es un valor de una variable aleatoria cuya distribución muestral es aproximadamente la distribución normal estándar.

- a) Con referencia al ejemplo 10.4, a un nivel de significancia de 0.05, pruebe si $\sigma = 10.0$ años contra la alternativa de que $\sigma \neq 10.0$ años.
- b) Se sospecha que la varianza de la distribución de las distancias en kilómetros que un modelo nuevo de automóvil equipado con un motor diesel recorre con 5 litros de combustible es menor que la varianza de la distribución de distancias que recorre el mismo modelo equipado con un motor de gasolina de 6 cilindros, la cual se sabe es $\sigma^2 = 6.25$. Si 72 recorridos de prueba con el modelo diesel tienen una varianza de 4.41, ¿podemos concluir, a un nivel de significancia de 0.05, que la varianza de las distancias recorridas por el modelo que funciona con diesel es menor que la del modelo que funciona con gasolina?

10.73 Se realiza un estudio para comparar el tiempo que les toma a hombres y mujeres ensamblar cierto producto. La experiencia indica que la distribución del tiempo tanto para hombres como para mujeres es aproximadamente normal, pero que la varianza del tiempo para las mujeres es menor que para los hombres. Una muestra aleatoria de los tiempos de 11 hombres y 14 mujeres produce los siguientes datos:

Hombres	Mujeres
$n_1 = 11$	$n_2 = 14$
$s_1 = 6.1$	$s_2 = 5.3$

Pruebe la hipótesis de que $\sigma_1^2 = \sigma_2^2$ contra la alternativa de que $\sigma_1^2 > \sigma_2^2$. Utilice un valor P en su conclusión.

10.74 En el ejercicio 10.41 de la página 358 pruebe la hipótesis a un nivel de significancia de 0.05 de que $\sigma_1^2 = \sigma_2^2$ contra la alternativa de que $\sigma_1^2 \neq \sigma_2^2$, donde σ_1^2 y σ_2^2 son las varianzas para el número de organismos por metro cuadrado de agua en los dos lugares diferentes de Cedar Run.

10.75 Remítase al ejercicio 10.39 de la página 358 y pruebe la hipótesis de que $\sigma_1^2 = \sigma_2^2$ contra la alternativa de que $\sigma_1^2 \neq \sigma_2^2$, donde σ_1^2 y σ_2^2 son las varianzas para la duración de las películas producidas por la empresa 1 y la empresa 2, respectivamente. Utilice un valor P .

10.76 Se comparan dos tipos de instrumentos para medir la cantidad de monóxido de azufre en la atmósfera en un experimento sobre la contaminación del

aire. Los investigadores desean determinar si los dos tipos de instrumentos proporcionan mediciones con la misma variabilidad. Se registran las siguientes lecturas para los dos instrumentos:

Monóxido de azufre	
Instrumento A	Instrumento B
0.86	0.87
0.82	0.74
0.75	0.63
0.61	0.55
0.89	0.76
0.64	0.70
0.81	0.69
0.68	0.57
0.65	0.53

Suponga que las poblaciones de mediciones se distribuyen de forma aproximadamente normal y pruebe la hipótesis de que $\sigma_A = \sigma_B$ contra la alternativa de que $\sigma_A \neq \sigma_B$. Use un valor P .

10.77 Se lleva a cabo un experimento para comparar el contenido de alcohol en una salsa de soya en dos líneas de producción diferentes. La producción se supervisa ocho veces al día. A continuación se presentan los datos.

Línea de producción 1.

0.48 0.39 0.42 0.52 0.40 0.48 0.52 0.52

Línea de producción 2.

0.38 0.37 0.39 0.41 0.38 0.39 0.40 0.39

Suponga que ambas poblaciones son normales. Se sospecha que la línea de producción 1 no está produciendo tan consistentemente como la línea 2 en términos de contenido de alcohol. Pruebe la hipótesis de que $\sigma_1 = \sigma_2$ contra la alternativa de que $\sigma_1 \neq \sigma_2$. Utilice un valor P .

10.78 Se sabe que las emisiones de hidrocarburos de los automóviles disminuyeron de forma drástica durante la década de 1980. Se realizó un estudio para comparar las emisiones de hidrocarburos a velocidad estacionaria, en partes por millón (ppm), para automóviles de 1980 y 1990. Se seleccionaron al azar 20 automóviles de cada modelo y se registraron sus niveles de emisión de hidrocarburos. Los datos son los siguientes:

Modelos 1980:

141 359 247 940 882 494 306 210 105 880
200 223 188 940 241 190 300 435 241 380

Modelos 1990:

140 160 20 20 223 60 20 95 360 70
220 400 217 58 235 380 200 175 85 65

Pruebe la hipótesis de que $\sigma_1 = \sigma_2$ contra la alternativa de que $\sigma_1 \neq \sigma_2$. Suponga que ambas poblaciones son normales. Utilice un valor P .

10.11 Prueba de la bondad de ajuste

A lo largo de este capítulo nos ocupamos de la prueba de hipótesis estadística acerca de parámetros de una sola población, como μ , σ^2 y p . Ahora consideraremos una prueba para determinar si una población tiene una distribución teórica específica. La prueba se basa en el nivel de ajuste que existe entre la frecuencia de ocurrencia de las observaciones en una muestra observada y las frecuencias esperadas que se obtienen a partir de la distribución hipotética.

Para ilustrar lo anterior considere el lanzamiento de un dado. Suponemos que se trata de un dado legal, lo cual equivale a probar la hipótesis de que la distribución de resultados es la distribución uniforme discreta

$$f(x) = \frac{1}{6}, \quad x = 1, 2, \dots, 6.$$

Suponga que el dado se lanza 120 veces y que se registra cada resultado. Teóricamente, si el dado está balanceado, esperaríamos que cada cara ocurriera 20 veces. Los resultados se presentan en la tabla 10.4.

Tabla 10.4: Frecuencias observadas y esperadas de 120 lanzamientos de un dado

Cara	1	2	3	4	5	6
Observadas	20	22	17	18	19	24
Esperadas	20	20	20	20	20	20

Al comparar las frecuencias observadas con las frecuencias esperadas correspondientes debemos decidir si es posible que tales discrepancias ocurran como resultado de fluctuaciones del muestreo, de que el dado está balanceado o no es legal o de que la distribución de resultados no es uniforme. Es práctica común referirse a cada resultado posible de un experimento como una celda. En nuestro caso tenemos 6 celdas. A continuación se define el estadístico adecuado en el cual basamos nuestro criterio de decisión para un experimento que incluye k celdas.

Una **prueba de la bondad de ajuste** entre las frecuencias observadas y esperadas se basa en la cantidad.

Prueba de la
bondad de
ajuste

$$\chi^2 = \sum_{i=1}^k \frac{(o_i - e_i)^2}{e_i},$$

donde χ^2 es un valor de una variable aleatoria cuya distribución muestral se aproxima muy de cerca a la distribución chi cuadrada con $v = k - 1$ grados de libertad. Los símbolos o_i y e_i representan las frecuencias observada y esperada, respectivamente, para la i -ésima celda.

El número de grados de libertad asociado con la distribución chi cuadrada que se utiliza aquí es igual a $k - 1$, pues sólo hay $k - 1$ frecuencias de celdas libremente determinadas. Es decir, una vez que se determinan las frecuencias de $k - 1$ celdas, también se determina la frecuencia para la k -ésima celda.

Si las frecuencias observadas se acercan a las frecuencias esperadas correspondientes, el valor χ^2 será pequeño, lo cual indica un buen ajuste. Si las frecuencias observadas difieren de manera considerable de las frecuencias esperadas, el valor χ^2 será grande y el ajuste deficiente. Un buen ajuste conduce a la aceptación de H_0 , mientras que un ajuste

deficiente conduce a su rechazo. Por lo tanto, la región crítica caerá en la cola derecha de la distribución chi cuadrada. Para un nivel de significancia igual a α encontramos el valor crítico χ^2_{α} de la tabla A.5 y, entonces, $\chi^2 > \chi^2_{\alpha}$ constituye la región crítica. El criterio de decisión que aquí se describe no se debería utilizar a menos que cada una de las frecuencias esperadas sea por lo menos igual a 5. Esta restricción podría requerir la combinación de celdas adyacentes, lo que dará como resultado una reducción en el número de grados de libertad.

En la tabla 10.4 encontramos que el valor χ^2 es

$$\chi^2 = \frac{(20-20)^2}{20} + \frac{(22-20)^2}{20} + \frac{(17-20)^2}{20} + \frac{(18-20)^2}{20} + \frac{(19-20)^2}{20} + \frac{(24-20)^2}{20} = 1.7.$$

Si usamos la tabla A.5, encontramos $\chi^2_{0.05} = 11.070$ para $v = 5$ grados de libertad. Como 1.7 es menor que el valor crítico, no se rechaza H_0 . Concluimos que no hay suficiente evidencia de que el dado está desbalanceado.

Como un segundo ejemplo probemos la hipótesis de que la distribución de frecuencias de la duración de baterías presentadas en la tabla 1.7 de la página 23 se puede aproximar mediante una distribución normal con media $\mu = 3.5$ y desviación estándar $\sigma = 0.7$. Las frecuencias esperadas para las 7 clases (celdas) que se listan en la tabla 10.5 se obtienen calculando las áreas bajo la curva normal hipotética que caen entre los diversos límites de clase.

Tabla 10.5: Frecuencias observadas y esperadas para la duración de las baterías suponiendo normalidad

Límites de clase	o_i	e_i
1.45-1.95	2	0.5
1.95-2.45	1	2.1
2.45-2.95	4	5.9
2.95-3.45	15	10.3
3.45-3.95	10	10.7
3.95-4.45	5	7.0
4.45-4.95	3	3.5

Por ejemplo, los valores z que corresponden a los límites de la cuarta clase son

$$z_1 = \frac{2.95 - 3.5}{0.7} = -0.79 \quad \text{y} \quad z_2 = \frac{3.45 - 3.5}{0.7} = -0.07.$$

En la tabla A.3 encontramos que el área entre $z_1 = -0.79$ y $z_2 = -0.07$ es

$$\begin{aligned} \text{área} &= P(-0.79 < Z < -0.07) = P(Z < -0.07) - P(Z < -0.79) \\ &= 0.4721 - 0.2148 = 0.2573. \end{aligned}$$

Por lo tanto, la frecuencia esperada para la cuarta clase es

$$e_4 = (0.2573)(40) = 10.3.$$

Se acostumbra redondear estas frecuencias a un decimal.

La frecuencia esperada para el primer intervalo de clase se obtiene utilizando el área total bajo la curva normal a la izquierda del límite 1.95. Para el último intervalo de clase usamos el área total a la derecha del límite 4.45. Todas las demás frecuencias esperadas se determinan utilizando el método que se describe para la cuarta clase. Observe que combinamos clases adyacentes en la tabla 10.5 donde las frecuencias esperadas son menores que 5 (una regla general en la prueba de la bondad de ajuste). En consecuencia, el número total de intervalos se reduce de 7 a 4, lo cual da como resultado $v = 3$ grados de libertad. Entonces, el valor χ^2 es dado por

$$\chi^2 = \frac{(7 - 8.5)^2}{8.5} + \frac{(15 - 10.3)^2}{10.3} + \frac{(10 - 10.7)^2}{10.7} + \frac{(8 - 10.5)^2}{10.5} = 3.05.$$

Como el valor χ^2 calculado es menor que $\chi_{0.05}^2 = 7.815$ para 3 grados de libertad, no tenemos razón para rechazar la hipótesis nula y concluimos que la distribución normal con $\mu = 3.5$ y $\sigma = 0.7$ proporciona un buen ajuste para la distribución de la duración de las baterías.

La prueba de bondad de ajuste chi cuadrada es un recurso importante, en particular debido a que muchos procedimientos estadísticos en la práctica dependen, en un sentido teórico, de la suposición de que los datos reunidos provienen de un tipo de distribución específico. Como ya se expuso, la suposición de normalidad se hace muy a menudo. En los siguientes capítulos continuaremos haciendo suposiciones de normalidad con el fin de proporcionar una base teórica para ciertas pruebas e intervalos de confianza.

En la literatura hay pruebas para evaluar la normalidad que son más poderosas que la prueba chi cuadrada. Una de tales pruebas es la **prueba de Geary**, la cual se basa en un estadístico muy sencillo que es el cociente de dos estimadores de la desviación estándar de la población σ . Suponga que se toma una muestra aleatoria X_1, X_2, \dots, X_n de una distribución normal, $N(\mu, \sigma)$. Considere el cociente

$$U = \frac{\sqrt{\pi/2} \sum_{i=1}^n |X_i - \bar{X}|/n}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2/n}}.$$

El lector debería reconocer que el denominador es un estimador razonable de σ sin importar si la distribución es normal o no. El numerador es un buen estimador de σ si la distribución es normal, pero podría sobrestimar o subestimar a σ cuando haya desviaciones de la normalidad. Así, los valores de U que difieren considerablemente de 1.0 representan la señal de que se debe rechazar la hipótesis de normalidad.

Para muestras grandes una prueba razonable se basa en la normalidad aproximada de U . El estadístico de prueba es, entonces, una estandarización de U dada por

$$Z = \frac{U - 1}{0.2661/\sqrt{n}}.$$

Desde luego, el procedimiento de prueba incluye la región crítica bilateral. Calculamos un valor de z a partir de los datos y no rechazamos la hipótesis de normalidad cuando

$$-z_{\alpha/2} < Z < z_{\alpha/2}.$$

En la bibliografía se cita un artículo que trata sobre la prueba de Geary (Geary, 1947).

10.12 Prueba de independencia (datos categóricos)

El procedimiento de prueba de chi cuadrada que se presentó en la sección 10.11 también se puede usar para probar la hipótesis de independencia de dos variables de clasificación. Suponga que deseamos determinar si las opiniones de los votantes residentes del estado de Illinois respecto a una nueva reforma fiscal son independientes de sus niveles de ingreso. Los sujetos de una muestra aleatoria de 1000 votantes registrados del estado de Illinois se clasifican de acuerdo con su posición en las categorías de ingreso bajo, medio o alto, y si están a favor o no de la nueva reforma fiscal. Las frecuencias observadas se presentan en la tabla 10.6, la cual se conoce como **tabla de contingencia**.

Tabla 10.6: Tabla de contingencia 2×3

Reforma fiscal	Nivel de ingreso			Total
	Bajo	Medio	Alto	
A favor	182	213	203	598
En contra	154	138	110	402
Total	336	351	313	1000

Una tabla de contingencia con r renglones y c columnas se denomina tabla $r \times c$ (" $r \times c$ " se lee " r por c "). Los totales de renglones y columnas en la tabla 10.6 se denominan **frecuencias marginales**. Nuestra decisión de aceptar o rechazar la hipótesis nula, H_0 , de que la opinión de un votante respecto a la nueva reforma fiscal es independiente de su nivel de ingreso, se basa en qué tan bien se ajusten las frecuencias observadas en cada una de las 6 celdas de la tabla 10.6 y en las frecuencias que esperaríamos para cada celda si supusiéramos que H_0 es verdadera. Para encontrar estas frecuencias esperadas definamos los siguientes eventos:

- L : Una persona seleccionada está en el nivel de ingresos bajo.
- M : Una persona seleccionada está en el nivel de ingresos medio.
- H : Una persona seleccionada está en el nivel de ingresos alto.
- F : Una persona seleccionada está a favor de la nueva reforma fiscal.
- A : Una persona seleccionada está en contra de la nueva reforma fiscal.

Podemos usar las frecuencias marginales para listar las siguientes estimaciones de probabilidad:

$$P(L) = \frac{336}{1000}, \quad P(M) = \frac{351}{1000}, \quad P(H) = \frac{313}{1000},$$

$$P(F) = \frac{598}{1000}, \quad P(A) = \frac{402}{1000}.$$

Ahora bien, si H_0 es verdadera y las dos variables son independientes, deberíamos tener

$$P(L \cap F) = P(L)P(F) = \left(\frac{336}{1000}\right)\left(\frac{598}{1000}\right),$$

$$P(L \cap A) = P(L)P(A) = \left(\frac{336}{1000}\right)\left(\frac{402}{1000}\right).$$

$$P(M \cap F) = P(M) P(F) = \left(\frac{351}{1000}\right) \left(\frac{598}{1000}\right),$$

$$P(M \cap A) = P(M) P(A) = \left(\frac{351}{1000}\right) \left(\frac{402}{1000}\right),$$

$$P(H \cap F) = P(H) P(F) = \left(\frac{313}{1000}\right) \left(\frac{598}{1000}\right),$$

$$P(H \cap A) = P(H) P(A) = \left(\frac{313}{1000}\right) \left(\frac{402}{1000}\right).$$

Las frecuencias esperadas se obtienen multiplicando la probabilidad de cada celda por el número total de observaciones. Como antes, redondeamos estas frecuencias a un decimal. Así, se estima que el número esperado de votantes de bajo ingreso en nuestra muestra que favorecen la reforma fiscal es

$$\left(\frac{336}{1000}\right) \left(\frac{598}{1000}\right) (1000) = \frac{(336)(598)}{1000} = 200.9$$

cuando H_0 es verdadera. La regla general para obtener la frecuencia esperada de cualquier celda es dada por la siguiente fórmula:

$$\text{frecuencia esperada} = \frac{(\text{total por columna}) \times (\text{total por renglón})}{\text{gran total}}$$

En la tabla 10.7 la frecuencia esperada para cada celda se registra entre paréntesis, a un lado del valor observado verdadero. Observe que las frecuencias esperadas en cualquier renglón o columna se suman al total marginal apropiado. En nuestro ejemplo necesitamos calcular sólo las dos frecuencias esperadas en el renglón superior de la tabla 10.7 y luego calcular las otras mediante sustracción. El número de grados de libertad asociados con la prueba chi cuadrada que aquí se usa es igual al número de frecuencias de celdas que se pueden llenar libremente cuando se nos proporcionan los totales marginales y el gran total, y en este caso ese número es 2. Una fórmula sencilla que proporciona el número correcto de grados de libertad es

$$v = (r - 1)(c - 1).$$

Tabla 10.7: Frecuencias observadas y esperadas

Reforma fiscal	Nivel de ingreso			Total
	Bajo	Medio	Alto	
A favor	182 (200.9)	213 (209.9)	203 (187.2)	598
En contra	154 (135.1)	138 (141.1)	110 (125.8)	402
Total	336	351	313	1000

Por lo tanto, para nuestro ejemplo $v = (2 - 1)(3 - 1) = 2$ grados de libertad. Para probar la hipótesis nula de independencia usamos el siguiente criterio de decisión:

Prueba de independencia Calcule

$$\chi^2 = \sum_i \frac{(o_i - e_i)^2}{e_i},$$

donde la sumatoria se extiende a todas las celdas rc en la tabla de contingencia $r \times c$.

Si $\chi^2 > \chi_{\alpha}^2$ con $\nu = (r-1)(c-1)$ grados de libertad, rechace la hipótesis nula de independencia al nivel de significancia α ; en otro caso no la rechace.

Al aplicar este criterio a nuestro ejemplo encontramos que

$$\begin{aligned} \chi^2 &= \frac{(182 - 200.9)^2}{200.9} + \frac{(213 - 209.9)^2}{209.9} + \frac{(203 - 187.2)^2}{187.2} \\ &\quad + \frac{(154 - 135.1)^2}{135.1} + \frac{(138 - 141.1)^2}{141.1} + \frac{(110 - 125.8)^2}{125.8} = 7.85, \\ P &\approx 0.02. \end{aligned}$$

En la tabla A.5 encontramos que $\chi_{0.05}^2 = 5.991$ para $\nu = (2-1)(3-1) = 2$ grados de libertad. Rechazamos la hipótesis nula y concluimos que la opinión de un votante respecto a la reforma fiscal y su nivel de ingresos no son independientes.

Es importante recordar que el estadístico sobre el cual basamos nuestra decisión tiene una distribución que sólo se aproxima por la distribución chi cuadrada. Los valores χ^2 calculados dependen de las frecuencias de las celdas y, en consecuencia, son discretos. La distribución chi cuadrada continua parece aproximarse muy bien a la distribución de muestreo discreta de χ^2 , siempre y cuando el número de grados de libertad sea mayor que 1. En una tabla de contingencia de 2×2 , donde sólo tenemos 1 grado de libertad, se aplica una corrección llamada **corrección de Yates para continuidad**.

La fórmula corregida entonces se convierte en

$$\chi^2(\text{corregida}) = \sum_i \frac{(|o_i - e_i| - 0.5)^2}{e_i}.$$

Si las frecuencias de las celdas esperadas son grandes, los resultados corregidos y sin corrección son casi iguales. Cuando las frecuencias esperadas están entre 5 y 10, se debe aplicar la corrección de Yates. Para frecuencias esperadas menores que 5 se debería utilizar la prueba exacta de Fisher-Irwin. Un análisis de esta prueba se puede encontrar en *Basic Concepts of Probability and Statistics* de Hodges y Lehmann (2005; véase la bibliografía). Sin embargo, la prueba de Fisher-Irwin se puede evitar seleccionando una muestra grande.

10.13 Prueba de homogeneidad

Cuando probamos la independencia en la sección 10.12 seleccionamos una muestra aleatoria de 1000 votantes, y determinamos al azar los totales de renglón y de columna para nuestra tabla de contingencia. Otro tipo de problema para el que se aplica el método de la sección 10.12 es aquel en el cual los totales de renglón y de columna están predeterminados. Suponga, por ejemplo, que decidimos de antemano seleccionar 200 demócratas, 150 republicanos y 150 independientes entre los votantes del estado de Carolina del Norte y registrar si están a favor de una iniciativa de ley para el aborto, si están en contra o si están indecisos. Las respuestas observadas se incluyen en la tabla 10.8.

Tabla 10.8: Frecuencias observadas

Ley para el aborto	Afiliación política			Total
	Demócrata	Republicano	Independiente	
A favor	82	70	62	214
En contra	93	62	67	222
Indeciso	25	18	21	64
Total	200	150	150	500

Ahora bien, en vez de hacer una prueba de independencia, probamos la hipótesis de que las proporciones de población dentro de cada renglón son iguales. Es decir, probamos la hipótesis de que las proporciones de demócratas, republicanos e independientes que están a favor de la ley para el aborto son iguales; las proporciones de cada afiliación política contra la ley son iguales y las proporciones de cada afiliación política que están indecisos son iguales. Básicamente nos interesamos en determinar si las tres categorías de votantes son **homogéneas** en lo que se refiere a sus opiniones acerca de la iniciativa de ley para el aborto. A esta prueba se le conoce como prueba de homogeneidad.

Al suponer homogeneidad de nuevo calculamos las frecuencias esperadas de las celdas multiplicando los totales de renglón y de columna correspondientes y después dividiendo entre el gran total. Luego continuamos el análisis utilizando el mismo estadístico chi cuadrada como antes. Ilustramos este proceso en el siguiente ejemplo para los datos de la tabla 10.8.

Ejemplo 10.14: Con respecto a los datos de la tabla 10.8 pruebe la hipótesis de que las opiniones en cuanto a la propuesta de ley para el aborto son las mismas en cada afiliación política. Utilice un nivel de significancia de 0.05.

- Solución:**
- H_0 : Para cada opinión las proporciones de demócratas, republicanos e independientes son iguales.
 - H_1 : Para al menos una opinión las proporciones de demócratas, republicanos e independientes no son iguales.
 - $\alpha = 0.05$.
 - Región crítica: $\chi^2 > 9.488$ con $v = 4$ grados de libertad.
 - Cálculos: necesitamos calcular las 4 frecuencias de las celdas usando la fórmula de las frecuencias de las celdas esperadas de la página 375. Todas las demás frecuencias se obtienen mediante sustracción. Las frecuencias de las celdas observadas y esperadas se muestran en la tabla 10.9.

Tabla 10.9: Frecuencias observadas y esperadas

Ley para el aborto	Afiliación política			Total
	Demócrata	Republicano	Independiente	
A favor	82 (85.6)	70 (64.2)	62 (64.2)	214
En contra	93 (88.8)	62 (66.6)	67 (66.6)	222
Indeciso	25 (25.6)	18 (19.2)	21 (19.2)	64
Total	200	150	150	500

Así,

$$\begin{aligned}\chi^2 &= \frac{(82 - 85.6)^2}{85.6} + \frac{(70 - 64.2)^2}{64.2} + \frac{(62 - 64.2)^2}{64.2} \\ &\quad + \frac{(93 - 88.8)^2}{88.8} + \frac{(62 - 66.6)^2}{66.6} + \frac{(67 - 66.6)^2}{66.6} \\ &\quad + \frac{(25 - 25.6)^2}{25.6} + \frac{(18 - 19.2)^2}{19.2} + \frac{(21 - 19.2)^2}{19.2} \\ &= 1.53.\end{aligned}$$

6. Decisión: No rechazar H_0 . No hay suficiente evidencia para concluir que la proporción de demócratas, republicanos e independientes difiere para cada opinión expresada. ▮

Prueba para varias proporciones

El estadístico chi cuadrada para probar la homogeneidad también se puede aplicar cuando se prueba la hipótesis de que k parámetros binomiales tienen el mismo valor. Por lo tanto, se trata de una extensión de la prueba que se presentó en la sección 10.9 para determinar las diferencias entre dos proporciones a una prueba para determinar diferencias entre k proporciones. En consecuencia, nos interesamos en probar la hipótesis nula

$$H_0: p_1 = p_2 = \dots = p_k$$

contra la hipótesis alternativa H_1 de que las proporciones de la población *no son todas iguales*. Para ejecutar esta prueba primero observamos muestras aleatorias independientes de tamaños n_1, n_2, \dots, n_k de las k poblaciones y ordenamos los datos en una tabla de contingencia $2 \times k$, la tabla 10.10.

Tabla 10.10: k muestras binomiales independientes

Muestra:	1	2	...	k
Éxitos	x_1	x_2	...	x_k
Fracasos	$n_1 - x_1$	$n_2 - x_2$...	$n_k - x_k$

De acuerdo con si los tamaños de las muestras aleatorias fueron predeterminados o si ocurrieron al azar, el procedimiento de prueba es idéntico a la prueba de homogeneidad o a la prueba de independencia. Por lo tanto, las frecuencias de las celdas esperadas se calculan como antes y se sustituyen junto con las frecuencias observadas en el estadístico chi cuadrada

$$\chi^2 = \sum_i \frac{(o_i - e_i)^2}{e_i},$$

con

$$v = (2 - 1)(k - 1) = k - 1$$

grados de libertad.

Al seleccionar la región crítica apropiada de la cola superior de la forma $\chi^2 > \chi^2_{\alpha}$ podemos llegar ahora a una decisión respecto a H_0 .

Ejemplo 10.15: En un estudio sobre un taller se reúne un conjunto de datos para determinar si la proporción de artículos defectuosos producida por los trabajadores fue la misma para el turno matutino, el vespertino y el nocturno. Los datos que se reunieron se muestran en la tabla 10.11.

Tabla 10.11: Datos para el ejemplo 10.15

Turno:	Matutino	Vespertino	Nocturno
Defectuosos	45	55	70
No defectuosos	905	890	870

Utilice un nivel de significancia de 0.025 para determinar si la proporción de artículos defectuosos es la misma para los tres turnos.

Solución: Representemos con p_1, p_2 y p_3 la proporción verdadera de artículos defectuosos para los turnos matutino, vespertino y nocturno, respectivamente.

1. $H_0: p_1 = p_2 = p_3$.
2. $H_1: p_1, p_2$ y p_3 no son iguales
3. $\alpha = 0.025$.
4. Región crítica: $\chi^2 > 7.378$ para $v = 2$ grados de libertad.
5. Cálculos: En correspondencia con las frecuencias observadas $o_1 = 45$ y $o_2 = 55$, encontramos

$$e_1 = \frac{(950)(170)}{2835} = 57.0 \quad \text{y} \quad e_2 = \frac{(945)(170)}{2835} = 56.7.$$

Todas las demás frecuencias esperadas se calculan restando y se incluyen en la tabla 10.12.

Tabla 10.12: Frecuencias esperadas y observadas

Turno:	Matutino	Vespertino	Nocturno	Total
Defectuosos	45 (57.0)	55 (56.7)	70 (56.3)	170
No defectuosos	905 (893.0)	890 (888.3)	870 (883.7)	2665
Total	950	945	940	2835

Ahora bien,

$$\begin{aligned} \chi^2 = & \frac{(45 - 57.0)^2}{57.0} + \frac{(55 - 56.7)^2}{56.7} + \frac{(70 - 56.3)^2}{56.3} \\ & + \frac{(905 - 893.0)^2}{893.0} + \frac{(890 - 888.3)^2}{888.3} + \frac{(870 - 883.7)^2}{883.7} = 6.29, \end{aligned}$$

$$P \approx 0.04.$$

6. Decisión: no rechazamos H_0 con $\alpha = 0.025$. Sin embargo, con el valor P calculado ciertamente sería riesgoso concluir que la proporción de artículos defectuosos producidos es la misma para todos los turnos. ■

A menudo un estudio completo implica utilizar métodos estadísticos en la prueba de hipótesis, lo que se puede mostrar a los ingenieros o científicos utilizando los

dos estadísticos de prueba, junto con valores P y gráficas estadísticas. Las gráficas complementan los diagnósticos numéricos con imágenes que indican de forma intuitiva por qué resultan esos valores P , así como qué tan razonables (o no) son las suposiciones operativas.

10.14 Estudio de caso de dos muestras

En esta sección consideramos un estudio que incluye un análisis gráfico y formal detallado, junto con la impresión por computadora con comentarios y conclusiones. En un estudio del análisis de datos que realizó el personal del Centro de Consulta Estadística del Virginia Tech se compararon dos materiales diferentes, la aleación A y la aleación B , en términos de la resistencia a la rotura. La aleación B es más costosa, aunque realmente se debería adoptar si se demuestra que es más fuerte que la aleación A . También se debe tomar en cuenta la consistencia del rendimiento de las dos aleaciones.

Se seleccionaron muestras aleatorias de vigas hechas con cada aleación y la resistencia se midió en unidades de flexión de 0.001 pulgadas cuando se aplicó una fuerza fija en ambos extremos de la viga. Se utilizaron 20 especímenes para cada una de las dos aleaciones. Los datos se presentan en la tabla 10.13.

Tabla 10.13: Datos para el estudio de caso de dos muestras

Aleación A			Aleación B		
88	82	87	75	81	80
79	85	90	77	78	81
84	88	83	86	78	77
89	80	81	84	82	78
81	85		80	80	
83	87		78	76	
82	80		83	85	
79	78		76	79	

Es importante que el ingeniero compare las dos aleaciones. Los investigadores están interesados en la resistencia y la reproducibilidad promedio, así como en determinar si hay una violación grave de la suposición de normalidad que requieren las pruebas t y F . Las figuras 10.21 y 10.22 son gráficas de cuantil-cuantil normales de las muestras de las dos aleaciones.

Al parecer no hay ninguna violación grave de la suposición de normalidad. Además, la figura 10.23 presenta dos gráficos de caja y bigote en la misma gráfica. Los gráficos de caja y bigote sugieren que no hay una diferencia apreciable en la variabilidad de la flexión para las dos aleaciones. Sin embargo, al parecer la flexión media de la aleación B es significativamente menor, lo cual sugiere (al menos gráficamente) que la aleación B es más fuerte. Las medias muestrales y las desviaciones estándar son

$$\bar{y}_A = 83.55, \quad s_A = 3.663; \quad \bar{y}_B = 79.70, \quad s_B = 3.097.$$

La impresión del SAS para el PROC TTEST se muestra en la figura 10.24. La prueba F sugiere que no hay una diferencia significativa en las varianzas ($P = 0.4709$) y el estadístico t de dos muestras para probar

$$H_0: \mu_A = \mu_B$$

$$H_1: \mu_A > \mu_B$$

($t = 3.59$, $P = 0.0009$) rechaza H_0 a favor de H_1 y, por consiguiente, confirma lo que sugiere la información gráfica. Aquí utilizamos la prueba t que agrupa las varianzas de dos muestras a la luz de los resultados de la prueba F . Con base en este análisis la adopción de la aleación B sería lo adecuado.

Significancia estadística y significancia científica o para la ingeniería

Mientras que el estadístico se podría sentir muy cómodo con los resultados de la comparación entre las dos aleaciones en el estudio de caso anterior, para el ingeniero queda un dilema. El análisis demostró una mejoría estadísticamente significativa utilizando la aleación B . Sin embargo, ¿realmente valdrá la pena aprovechar la diferencia que se en-

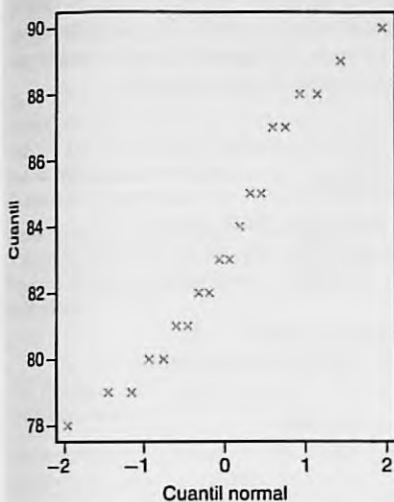


Figura 10.21: Gráfica de cuantil-cuantil normal de los datos para la aleación A .

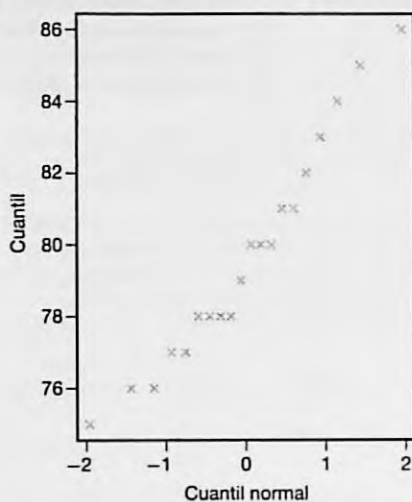


Figura 10.22: Gráfica de cuantil-cuantil normal de los datos para la aleación B .

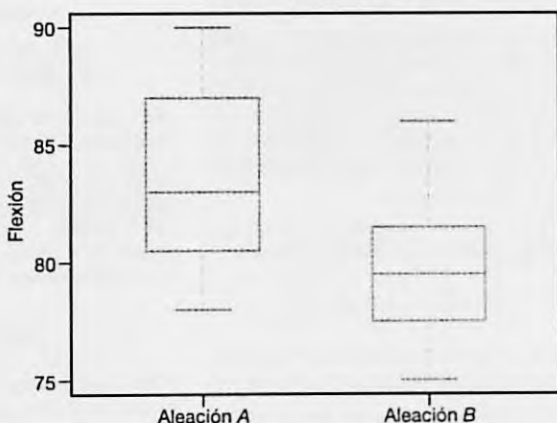


Figura 10.23: Gráficos de caja y bigote para ambas aleaciones.

contró si la aleación *B* es más costosa? Este ejemplo resalta una cuestión muy importante que con frecuencia pasan por alto los estadísticos y los analistas de datos: *la diferencia entre significancia estadística y significancia científica o para la ingeniería*. Aquí la diferencia promedio en la flexión es $\bar{y}_A - \bar{y}_B = 0.00385$ pulgadas. En un análisis completo el ingeniero debe determinar si la diferencia es suficiente para justificar el costo adicional a largo plazo. Ésta es una cuestión económica y de ingeniería. El lector debería comprender que una diferencia significativa en términos estadísticos tan sólo implica que la diferencia en las medias muestrales que se encuentra en los datos difícilmente podría ocurrir por casualidad. Esto no implica que la diferencia en las medias de la población sea profunda o particularmente significativa en el contexto del problema. Por ejemplo, en la sección 10.4 se utilizó una impresión por computadora con comentarios para demostrar la evidencia de que un medidor de pH está, de hecho, sesgado. Es decir, esto no demuestra un pH promedio de 7.00 para el material en que se probó. Pero la variabilidad entre las observaciones en la muestra es muy pequeña. El ingeniero podría decidir que las desviaciones pequeñas de 7.0 representan el medidor de pH adecuado.

The TTEST Procedure				
Alloy	N	Mean	Std Dev	Std Err
Alloy A	20	83.55	3.6631	0.8191
Alloy B	20	79.7	3.0967	0.6924
Variances	DF	t Value	Pr > t	
Equal	38	3.59	0.0009	
Unequal	37	3.59	0.0010	
Equality of Variances				
Num DF	Den DF	F Value	Pr > F	
19	19	1.40	0.4709	

Figura 10.24: Impresión del SAS con comentarios para los datos de las aleaciones.

Ejercicios

10.79 Se supone que una máquina mezcla cacahuates, avellanas, castañas y pacanas a razón de 5:2:2:1. Se observa que una lata que contiene 500 de tales nueces mezcladas tiene 269 cacahuates, 112 avellanas, 74 castañas y 45 pacanas. A un nivel de significancia de 0.05 pruebe la hipótesis de que la máquina mezcla las nueces a una razón de 5:2:2:1.

10.80 Las calificaciones de un curso de estadística para un semestre específico fueron las siguientes:

Calificación	A	B	C	D	F
<i>f</i>	14	18	32	20	16

Pruebe la hipótesis, a un nivel de significancia de 0.05, de que la distribución de calificaciones es uniforme.

10.81 Se lanza un dado 180 veces con los siguientes resultados:

<i>x</i>	1	2	3	4	5	6
<i>f</i>	28	36	36	30	27	23

¿Se trata de un dado balanceado? Utilice un nivel de significancia de 0.01.

10.82 Se seleccionan tres canicas de una urna que contiene 5 canicas rojas y 3 verdes. Después de registrar el número *X* de canicas rojas, las canicas se reemplazan en la urna y el experimento se repite 112 veces. Los resultados que se obtienen son los siguientes:

<i>x</i>	0	1	2	3
<i>f</i>	1	31	55	25

A un nivel de significancia de 0.05, pruebe la hipótesis de que los datos registrados se pueden ajustar a la distribución hipergeométrica $h(x; 8, 3, 5)$, $x = 0, 1, 2, 3$.

10.83 Se lanza una moneda hasta que sale una cara y se registra el número de lanzamientos X . Después de repetir el experimento 256 veces, obtenemos los siguientes resultados:

x	1	2	3	4	5	6	7	8
f	136	60	34	12	9	1	3	1

A un nivel de significancia de 0.05, pruebe la hipótesis de que la distribución observada de X se puede ajustar a la distribución geométrica $g(x; 1/2)$, $x = 1, 2, 3, \dots$

10.84 En el ejercicio 1.18 de la página 31 pruebe la bondad de ajuste entre las frecuencias de clase observadas y las frecuencias esperadas correspondientes de una distribución normal con $\mu = 65$ y $\sigma = 21$. Utilice un nivel de significancia de 0.05.

10.85 En el ejercicio 1.19 de la página 31 pruebe la bondad de ajuste entre las frecuencias de clase observadas y las frecuencias esperadas correspondientes de una distribución normal con $\mu = 1.8$ y $\sigma = 0.4$. Utilice un nivel de significancia de 0.01.

10.86 En un experimento diseñado para estudiar la dependencia de la hipertensión con respecto a los hábitos de fumar se tomaron los siguientes datos de 180 individuos:

	No fumadores	Fumadores moderados	Fumadores empedernidos
Con hipertensión	21	36	30
Sin hipertensión	48	26	19

Pruebe la hipótesis de que la presencia o ausencia de hipertensión es independiente de los hábitos de tabaquismo. Utilice un nivel de significancia de 0.05.

10.87 Una muestra aleatoria de 90 adultos se clasifica de acuerdo con el género y el número de horas dedicadas a ver la televisión durante una semana:

	Género	
	Masculino	Femenino
Más de 25 horas	15	29
Menos de 25 horas	27	19

Utilice un nivel de significancia de 0.01 y pruebe la hipótesis de que el tiempo dedicado a ver la televisión es independiente de si el espectador es hombre o mujer.

10.88 Una muestra aleatoria de 200 hombres casados, todos jubilados, se clasificó de acuerdo con la educación y el número de hijos:

Educación	Número de hijos		
	0-1	2-3	Más de 3
Primaria	14	37	32
Secundaria	19	42	17
Universidad	12	17	10

Utilice un nivel de significancia de 0.05 para probar la hipótesis de que el tamaño de la familia es independiente del nivel académico del padre.

10.89 Un criminólogo realizó una investigación para determinar si la incidencia de ciertos tipos de delitos varía de una parte de una gran ciudad a otra. Los crímenes específicos de interés eran el asalto, el robo de casas, el hurto y el homicidio. La siguiente tabla muestra el número de delitos cometidos en cuatro áreas de la ciudad durante el año pasado.

Distrito	Tipo de crimen			
	Asalto	Robo de casas	Hurto	Homicidio
1	162	118	451	18
2	310	196	996	25
3	258	193	458	10
4	280	175	390	19

¿A partir de estos datos podemos concluir, a un nivel de significancia de 0.01, que la ocurrencia de estos tipos de delitos depende del distrito de la ciudad?

10.90 De acuerdo con un estudio de la Universidad Johns Hopkins, publicado en *American Journal of Public Health*, las viudas viven más que los viudos. Considere los siguientes datos reunidos de supervivencia de 100 viudas y 100 viudos después de la muerte del cónyuge:

Años vividos	Viuda	Viudo
Menos de 5	25	39
de 5 a 10	42	40
Más de 10	33	21

Con un nivel de significancia de 0.05, ¿podemos concluir que las proporciones de viudas y viudos son iguales con respecto a los diferentes periodos que un cónyuge sobrevive luego de la muerte de su compañero?

10.91 Las siguientes respuestas respecto al nivel de vida en el momento en que se aplicó una encuesta de opinión independiente a 1000 familias, comparadas con sus respuestas sobre su nivel de vida del año anterior, parecen coincidir con los resultados de un estudio publicado en *Across the Board* (junio de 1981):

Periodo	Nivel de vida			Total	
	Un poco mejor	Igual	No tan bueno		
1980:	Ene.	72	144	84	300
	May	63	135	102	300
	Sept.	47	100	53	200
1981:	Ene.	40	105	55	200

Pruebe la hipótesis de que las proporciones de familias dentro de cada nivel de vida son iguales para cada uno de los cuatro periodos. Utilice un valor P .

10.92 La enfermería de una universidad realizó un experimento para determinar el grado de alivio que brindan tres jarabes para la tos. Cada jarabe se probó en 50 estudiantes y se registraron los siguientes datos:

	Jarabe para la tos		
	NyQuil	Robitussin	Triaminic
Sin alivio	11	13	9
Cierto alivio	32	28	27
Alivio completo	7	9	14

Pruebe la hipótesis de que los tres remedios para la tos son igualmente efectivos. Utilice un valor P en sus conclusiones.

10.93 Para determinar las posturas actuales acerca de rezar en escuelas públicas se llevó a cabo una investigación en 4 condados de Virginia. En la siguiente tabla se presentan las opiniones de 200 padres del condado de Craig, de 150 padres del condado de Giles, de 100 padres del condado de Franklin y de 100 padres del condado de Montgomery:

Actitud	Condado			
	Craig	Giles	Franklin	Mont.
A favor	65	66	40	34
En contra	42	30	33	42
Sin opinión	93	54	27	24

Pruebe la homogeneidad de las posturas entre los 4 condados respecto a rezar en escuelas públicas. Utilice un valor P en sus conclusiones.

10.94 Se lleva a cabo una encuesta en Indiana, Kentucky y Ohio para determinar la postura de los votantes respecto al transporte escolar. Un grupo de 200 votantes de cada uno de estos estados proporcionó los siguientes resultados:

Ejercicios de repaso

10.97 Plantee las hipótesis nula y alternativa que utilizaría para probar las siguientes afirmaciones y determine de manera general en dónde se localiza la región crítica:

- La cantidad promedio de nieve que cae en el lago George durante el mes de febrero es de 21.8 centímetros.
- No más del 20% de los profesores de la universidad local contribuyó al fondo anual para donaciones.
- En promedio, los niños asisten a la escuela en un área de 6.2 kilómetros de sus casas en un suburbio de St. Louis.
- Al menos 70% de los automóviles nuevos del siguiente año caerán en la categoría de compactos y semicompactos.
- La proporción de votantes que están a favor del

Estado	Postura del votante		
	Apoya	No apoya	Indeciso
Indiana	82	97	21
Kentucky	107	66	27
Ohio	93	74	33

A un nivel de significancia de 0.05 pruebe la hipótesis nula de que las proporciones de votantes dentro de cada categoría de postura son las mismas en cada uno de los tres estados.

10.95 Se lleva a cabo una investigación en dos ciudades de Virginia para determinar la opinión de los votantes respecto a dos candidatos a la gubernatura en una elección próxima. En cada ciudad se seleccionaron 500 votantes al azar y se registraron los siguientes datos:

Opinión del votante	Ciudad	
	Richmond	Norfolk
A favor de A	204	225
A favor de B	211	198
Indeciso	85	77

A un nivel de significancia de 0.05 pruebe la hipótesis nula de que las proporciones de votantes que están a favor del candidato A, a favor del candidato B o que están indecisos son las mismas para cada ciudad.

10.96 En un estudio para estimar la proporción de esposas que de manera regular ven telenovelas se encuentra que 52 de 200 esposas en Denver, 31 de 150 en Phoenix y 37 de 150 en Rochester ven al menos una telenovela. Utilice un nivel de significancia de 0.05 para probar la hipótesis de que no hay diferencia entre las proporciones verdaderas de esposas que ven telenovelas en esas tres ciudades.

funcionario actual para la próxima elección es de 0.58.

- El filete rib-eye promedio en el restaurante Longhorn Steak pesa al menos 340 gramos.

10.98 Un genetista se interesa en la proporción de hombres y mujeres de una población que tiene cierto trastorno sanguíneo menor. En una muestra aleatoria de 100 hombres se encuentra que 31 lo padecen, mientras que sólo 24 de 100 mujeres analizadas tienen el trastorno. Con un nivel de significancia de 0.01, ¿podemos concluir que la proporción de hombres en la población con este trastorno sanguíneo es significativamente mayor que la proporción de mujeres afectadas?

10.99 Se realizó un estudio para determinar si un número mayor de italianos que de estadounidenses prefieren la champaña blanca en vez de la rosa para

las bodas. De los 300 italianos que se seleccionaron al azar, 72 preferían champaña blanca, y de los 400 estadounidenses seleccionados, 70 preferían champaña blanca en vez de la rosa. ¿Podemos concluir que una proporción mayor de italianos que de estadounidenses prefiere champaña blanca en las bodas? Utilice un nivel de significancia de 0.05.

10.100 Considere la situación del ejercicio 10.54 de la página 360. También se midió el consumo de oxígeno en mL/kg/min.

Sujeto	Con CO	Sin CO
1	26.46	25.41
2	17.46	22.53
3	16.32	16.32
4	20.19	27.48
5	19.84	24.97
6	20.65	21.77
7	28.21	28.17
8	33.94	32.02
9	29.32	28.96

Se supone que el consumo de oxígeno debería ser mayor en un ambiente relativamente libre de CO. Realice una prueba de significancia y analice la suposición.

10.101 En un estudio realizado por el Centro de Consulta Estadística de Virginia Tech se solicitó a un grupo de sujetos realizar cierta tarea en la computadora. La respuesta que se midió fue el tiempo requerido para realizar la tarea. El propósito del experimento fue probar un grupo de herramientas de ayuda desarrolladas por el Departamento de Ciencias Computacionales de la universidad. En el estudio participaron 10 sujetos. Con una asignación al azar, a 5 se les dio un procedimiento estándar usando lenguaje Fortran para realizar la tarea. A los otros 5 se les pidió realizar la tarea usando las herramientas de ayuda. A continuación se presentan los datos del tiempo requerido para completar la tarea.

Grupo 1 (procedimiento estándar)	Grupo 2 (herramienta de ayuda)
161	132
169	162
174	134
158	138
163	133

Suponga que las distribuciones de la población son normales y las varianzas son las mismas para los dos grupos y apoye o refute la conjetura de que las herramientas de ayuda aumentan la velocidad con la que se realiza la tarea.

10.102 Establezca las hipótesis nula y alternativa que usaría para probar las siguientes afirmaciones, y determine de manera general en dónde se localiza la región crítica:

- A lo sumo, 20% de la cosecha de trigo del próximo año se exportará a la Unión Soviética.
- En promedio, las amas de casa estadounidenses beben 3 tazas de café al día.
- La proporción de estudiantes que se graduaron este año en Virginia, especializados en ciencias sociales, es de al menos 0.15.
- El donativo promedio a la American Lung Association no es mayor de 10 dólares.
- Los residentes de la zona suburbana de Richmond viajan en promedio 15 kilómetros para llegar a su lugar de trabajo.

10.103 Si se selecciona al azar una lata que contiene 500 nueces de cada uno de tres distribuidores de nueces surtidas y cada lata contiene 345, 313 y 359 cacahuates, respectivamente. Con un nivel de significancia de 0.01, ¿podríamos concluir que las nueces surtidas de los tres distribuidores contienen proporciones iguales de cacahuates?

10.104 Se realiza un estudio para determinar si hay una diferencia entre las proporciones de padres en los estados de Maryland (MD), Virginia (VA), Georgia (GA) y Alabama (AL) que están a favor de colocar Biblias en las escuelas primarias. En la siguiente tabla se registran las respuestas de 100 padres seleccionados al azar en cada uno de esos estados:

Preferencia	Estado			
	MD	VA	GA	AL
Sí	65	71	78	82
No	35	29	22	18

¿Podemos concluir que las proporciones de padres que están a favor de colocar Biblias en las escuelas son iguales en esos cuatro estados? Utilice un nivel de significancia de 0.01.

10.105 Se lleva a cabo un estudio en el Centro de Medicina Veterinaria Equina de la Universidad Regional de Virginia en Maryland para determinar si la realización de cierto tipo de cirugía en caballos jóvenes tiene algún efecto en ciertas clases de células sanguíneas del animal. Se toman muestras del fluido de seis potros antes y después de la cirugía. En las muestras se analiza el número de leucocitos de glóbulos blancos (GB) después de la operación. También se midieron los leucocitos GB preoperatorios. Los datos son los siguientes:

Potro	Precirugía*	Postcirugía*
1	10.80	10.60
2	12.90	16.60
3	9.59	17.20
4	8.81	14.00
5	12.00	10.60
6	6.07	8.60

*Todos los valores $\times 10^{-3}$

Utilice una prueba t de una muestra pareada para determinar si hay un cambio significativo en los leucocitos GB con la cirugía.

10.106 El Departamento de Salud y Educación Física de Virginia Tech realizó un estudio para determinar si 8 semanas de entrenamiento realmente reducen los niveles de colesterol de los participantes. A un grupo de tratamiento que consta de 15 personas se les dieron conferencias dos veces a la semana acerca de cómo reducir sus niveles de colesterol. Otro grupo de 18 personas, de edad similar, fue seleccionado al azar como grupo de control. Se registraron los siguientes niveles de colesterol de todos los participantes al final del programa de 8 semanas:

Grupo con tratamiento:

Tratamiento:

129 131 154 172 115 126 175 191
122 238 159 156 176 175 126

Control:

151 132 196 195 188 198 187 168 115
165 137 208 133 217 191 193 140 146

¿Podemos concluir, a un nivel de significancia del 5%, que el nivel de colesterol promedio se redujo gracias al programa? Haga la prueba adecuada en las medias.

10.107 En un estudio que llevó a cabo el Departamento de Ingeniería Mecánica, el cual fue analizado por el Centro de Consulta Estadística del Virginia Tech, se compararon las varillas de acero distribuidas por dos empresas diferentes. Se fabricaron diez resortes de muestra con las varillas proporcionadas por cada empresa y se estudió la "capacidad de rebote". Los datos son los siguientes:

Empresa A:

9.3 8.8 6.8 8.7 8.5 6.7 8.0 6.5 9.2 7.0

Empresa B:

11.0 9.8 9.9 10.2 10.1 9.7 11.0 11.1 10.2 9.6

¿Puede concluir que casi no hay diferencia en las medias entre las varillas de acero proporcionadas por las dos empresas? Utilice un valor P para llegar a su conclusión. ¿Deberían agruparse las varianzas en este caso?

10.108 En un estudio realizado por el Centro de Recursos Acuáticos, el cual fue analizado por el Centro de Consulta Estadística del Virginia Tech, se com-

pararon dos diferentes plantas de tratamiento para aguas residuales. La planta A se ubica en una zona donde el ingreso medio de los hogares está por abajo de \$22,000 al año, y la planta B se ubica en un lugar donde el ingreso medio de los hogares está por arriba de \$60,000 anuales. La cantidad de agua residual tratada en cada planta (miles de galones/día) se muestreó de forma aleatoria durante 10 días. Los datos son los siguientes:

Planta A:

21 19 20 23 22 28 32 19 13 18

Planta B:

20 39 24 33 30 28 30 22 33 24

A un nivel de significancia de 5%, ¿podemos concluir que la cantidad promedio de agua residual tratada en la planta del vecindario de altos ingresos es mayor que la tratada en la planta del área de bajos ingresos? Suponga normalidad.

10.109 Los siguientes datos muestran el número de defectos en 100,000 líneas de código en un tipo particular de software hecho en Estados Unidos y en Japón. ¿Hay suficiente evidencia para afirmar que existe una diferencia significativa entre los programas creados en los dos países? Pruebe las medias. ¿Se deberían agrupar las varianzas?

Estados Unidos	48	39	42	52	40	48	52	52
Japón	54	48	52	55	43	46	48	52
	50	48	42	40	43	48	50	46
	38	38	36	40	40	48	48	45

10.110 Existen estudios que muestran que la concentración de PCB es mucho más alta en tejido mamario maligno que en tejido mamario normal. Si un estudio de 50 mujeres con cáncer de mama revela una concentración promedio de PCB de 22.8×10^{-4} gramos, con una desviación estándar de 4.8×10^{-4} gramos, ¿la concentración media de PCB es menor que 24×10^{-4} gramos?

10.111 Valor z para probar $p_1 - p_2 = d_0$: Para probar la hipótesis nula H_0 de que $p_1 - p_2 = d_0$, donde $d_0 \neq 0$, basamos nuestra decisión en

$$z = \frac{\hat{p}_1 - \hat{p}_2 - d_0}{\sqrt{\hat{p}_1 \hat{q}_1 / n_1 + \hat{p}_2 \hat{q}_2 / n_2}}$$

que es un valor de una variable aleatoria cuya distribución se aproxima a la distribución normal estándar, siempre y cuando n_1 y n_2 sean grandes. Con respecto al ejemplo 10.11 de la página 364, pruebe la hipótesis de que el porcentaje de votantes de la ciudad que están a favor de la construcción de la planta química no excederá en más de 3% al porcentaje de votantes del condado. Utilice un valor P en su conclusión.

10.15 Posibles riesgos y errores conceptuales; relación con el material de otros capítulos

Una de las formas más sencillas de darle un uso incorrecto a la estadística se refiere a la conclusión científica final que se obtiene cuando el analista no rechaza la hipótesis nula H_0 . En este texto intentamos aclarar lo que significan la hipótesis nula y la alternativa, y también enfatizamos que, en general, la hipótesis alternativa es mucho más importante. A modo de ejemplo, si un ingeniero trata de comparar dos calibradores utilizando una prueba t de dos muestras, y H_0 afirma que "los calibradores son equivalentes", mientras que H_1 afirma que "los calibradores no son equivalentes", no rechazar H_0 no lo lleva a concluir que los calibradores son equivalentes. De hecho, ¡se puede dar el caso de que nunca se escriba o se diga "acepto H_0 "! El hecho de no rechazar H_0 sólo implica que no existe evidencia suficiente. Según la naturaleza de la hipótesis, no se descartan aún muchas posibilidades.

En el capítulo 9 consideramos el caso del intervalo de confianza para muestras grandes utilizando

$$z = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

En la prueba de hipótesis es riesgoso reemplazar σ con s para $n < 30$. Si $n \geq 30$ y la distribución no es normal pero se acerca hasta cierto punto a la normal, se requiere el teorema del límite central y se confía en el hecho de que con $n \geq 30$, $s \approx \sigma$. Desde luego, cualquier prueba t va acompañada por la suposición concomitante de normalidad. Como en el caso de los intervalos de confianza, la prueba t es relativamente robusta para la normalidad. Sin embargo, cuando la muestra no es demasiado pequeña es necesario utilizar gráficas de probabilidad normal, pruebas de bondad de ajuste u otros procedimientos gráficos.

La mayoría de los capítulos de este texto incluyen análisis que tienen el propósito de relacionar el capítulo en cuestión con el siguiente material. Los temas de estimación y prueba de hipótesis se utilizan de manera importante en casi todas las técnicas que entran en el concepto de "métodos estadísticos". Los estudiantes lo notarán fácilmente cuando avancen a los capítulos 11 a 16. Será evidente que esos capítulos dependen en gran medida de los modelos estadísticos. Los estudiantes se verán expuestos al uso de los modelos en una gran variedad de aplicaciones, en diversos campos científicos y de la ingeniería. Rápidamente se darán cuenta de que el esquema de un modelo estadístico es inútil a menos que se disponga de datos para estimar parámetros en el modelo formulado. Esto será especialmente evidente en los capítulos 11 y 12, cuando se presente el concepto de modelos de regresión. Seguiremos utilizando los conceptos y la teoría relacionados con el capítulo 9. En lo que se refiere al material de este capítulo, el esquema de la prueba de hipótesis, de los valores P , de la potencia de una prueba y la selección del tamaño de la muestra, en conjunto desempeñarán un papel importante. Dado que con mucha frecuencia la formulación del modelo inicial debe complementarse con la edición del mismo antes de que el analista se sienta lo suficientemente cómodo para utilizarlo con el fin de conocer o predecir un proceso, en los capítulos 11, 12 y 15 se utilizará con frecuencia la prueba de hipótesis para complementar las medidas diagnósticas que se emplean con el fin de evaluar la calidad del modelo.