

SIMULACIÓN Y LA MÉTODO MONTE CARLO

SERIE DE WILEY EN PROBABILIDAD Y ESTADÍSTICA

Establecido por *Walter A. Shewhart y Samuel S. Wilks*

Editores: *David J. Balding, Noel AC Cressie, Garrett M. Fitzmaurice, Geof H. Givens, Harvey Goldstein, Geert Molenberghs, David W. Scott, Adrian FM Smith, Ruey S. Tsay, Sanford Weisberg*

Editores eméritos: *J. Stuart Hunter, Iain M. Johnstone, Joseph B. Kadane, Jozef L. Teugels*

La Serie de Wiley en Probabilidad y Estadística está bien establecida y tiene autoridad. Cubre muchos temas de interés de investigación actual tanto en estadística pura como aplicada y teoría de la probabilidad. Escritos por destacados estadísticos e instituciones, los títulos abarcan tanto los desarrollos más avanzados en el campo como los métodos clásicos.

Como reflejo de la amplia gama de investigaciones actuales en estadística, la serie abarca estadísticas aplicadas, metodológicas y teóricas, que van desde aplicaciones y nuevas técnicas posibles gracias a los avances en la práctica computarizada hasta el tratamiento riguroso de los enfoques teóricos. Esta serie proporciona una lectura esencial e invaluable para todos los estadísticos, ya sea en la academia, la industria, el gobierno o la investigación.

Puede encontrar una lista completa de los títulos de esta serie en <http://www.wiley.com/go/wsps>

SIMULACIÓN Y LA MÉTODO MONTE CARLO

Tercera edición

Reuven Y Rubinstein

Technión

dirk p kroese

Universidad de Queensland

WILEY

Copyright © 2017 por John Wiley & Sons, Inc. Todos los derechos reservados.

Publicado por John Wiley & Sons, Inc., Hoboken, Nueva Jersey.

Publicado simultáneamente en Canadá.

Ninguna parte de esta publicación puede reproducirse, almacenarse en un sistema de recuperación o transmitirse de ninguna forma ni por ningún medio, ya sea electrónico, mecánico, fotocopiado, grabado, escaneado o cualquier otro, excepto según lo permita la Sección 107 o 108 de la Ley de Estados Unidos de 1976. Ley de derechos de autor de los Estados Unidos, sin el permiso previo por escrito del editor o la autorización mediante el pago de la tarifa correspondiente por copia a Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, o en la web en www.copyright.com. Las solicitudes de permiso al editor deben dirigirse al Departamento de permisos, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, o en línea en <http://www.wiley.com/go/permission>.

Límite de responsabilidad/Descargo de responsabilidad de la garantía: Si bien el editor y el autor han realizado sus mejores esfuerzos para preparar este libro, no hacen declaraciones ni garantías con respecto a la precisión o integridad del contenido de este libro y específicamente renuncian a cualquier garantía implícita de comerciabilidad, o aptitud para un propósito particular. Los representantes de ventas o los materiales de ventas por escrito no pueden crear ni extender ninguna garantía. Los consejos y estrategias contenidos en este documento pueden no ser adecuados para su situación. Deberías consultar con un profesional apropiado. Ni el editor ni el autor serán responsables de ninguna pérdida de ganancias ni de ningún otro daño comercial, incluidos, entre otros, daños especiales, incidentales, consecuentes o de otro tipo.

Para obtener información general sobre nuestros otros productos y servicios o para soporte técnico, comuníquese con nuestro Departamento de Atención al Cliente dentro de los Estados Unidos al (800) 762-2974, fuera de los Estados Unidos al (317) 572-3993 o fax (317) 572-4002.

Wiley también publica sus libros en una variedad de formatos electrónicos. Algunos contenidos que aparecen impresos pueden no estar disponibles en formato electrónico. Para obtener más información acerca de los productos Wiley, visite nuestro sitio web en www.wiley.com.

Datos de catalogación en publicación de la Biblioteca del Congreso:

Nombres: Rubinstein, Reuven Y. | Kroese, Dirk

P. Título: Simulación y el método Monte Carlo.

Descripción: Tercera edición / Reuven Rubinstein, Dirk P. Kroese. | hoboken, Nueva Jersey: John Wiley & Sons, Inc., [2017] | Serie: Serie de Wiley en probabilidad y estadística | Incluye referencias bibliográficas e índice. Identificadores: LCCN

2016020742 (imprimir) | LCCN 2016023293 (libro electrónico) | ISBN

9781118632161 (tela) | ISBN 9781118632208 (pdf) | ISBN 9781118632383 (epub)

Temas: LCSH: Método Monte Carlo. | Simulación por computadora digital. |

Estadística matemática. | Muestreo (Estadísticas)

Clasificación: LCC QA298 .R8 2017 (imprimir) | LCC QA298 (libro electrónico) | DDC 518/.282--dc23

Registro LC disponible en <https://lcn.loc.gov/2016020742>

Impreso en los Estados Unidos de América.

10 9 8 7 6 5 4 3 2 1

A la memoria de
Reuven Y. Rubinstein (1938-2012)



Reuven Rubinstein falleció durante la redacción de esta tercera edición. Reuven fue uno de los pioneros de la simulación de Monte Carlo y permaneció a la vanguardia de la investigación en esta área hasta el final de su vida. En 2011 recibió el más alto honor otorgado por INFORMS Simulation Society: The Lifetime Professional Achievement Award, donde sus logros se resumen de la siguiente manera:

El profesor Rubinstein ha sido una figura fundamental en la teoría y la práctica de la simulación tal como la conocemos hoy. Su carrera refleja un alto nivel de creatividad y contribución, con voluntad de explorar nuevas áreas y una asombrosa capacidad para sugerir nuevas vías de investigación sorprendentes e influir en el trabajo posterior.

Que su contagioso entusiasmo y curiosidad vivan a través de sus libros y de las muchas personas a las que inspiró.

dirk p kroese

CONTENIDO

Prefacio	XIII
Expresiones de gratitud	xvii
1 Preliminares	1
1.1 Introducción	1
1.2 Experimentos aleatorios	1
1.3 Probabilidad condicional e independencia	2
1.4 Variables aleatorias y distribuciones de probabilidad	4
1.5 Algunas distribuciones importantes	5
1.6 Expectativa	6
1.7 Distribuciones conjuntas	7
1.8 Funciones de Variables Aleatorias	11
1.8.1 Transformaciones lineales	12
1.8.2 Transformaciones generales	13
1.9 Transformaciones	14
1.10 Variables aleatorias conjuntas normales	15
1.11 Teoremas de los límites	dieciséis
1.12 Procesos de Veneno	17
1.13 Procesos de Markov	19
1.13.1 Cadenas de Markov	19
1.13.2 Clasificación de Estados	21
	viii

1.13.3 Comportamiento limitante	22
1.13.4 Reversibilidad	24
1.13.5 Procesos de salto de Markov	25
1.14 Procesos gaussianos	27
1.15 Información	28
1.15.1 Entropía de Shannon	29
1.15.2 Entropía cruzada de Kullback-Leibler	31
1.15.3 Estimador de máxima verosimilitud y función de puntuación	32
1.15.4 Información del pescador	33
1.16 Optimización Convexa y Dualidad	34
1.16.1 Método Lagrangiano	35
1.16.2 Dualidad	37
Problemas	41
Referencias	46
2 Número aleatorio, variable aleatoria y generación de procesos estocásticos	49
2.1 Introducción	49
2.2 Generación de números aleatorios	49
2.2.1 Generadores recursivos múltiples	51
2.2.2 Generadores lineales de módulo 2	52
2.3 Generación de variables aleatorias	55
2.3.1 Método de transformada inversa	55
2.3.2 Método de alias	57
2.3.3 Método de composición	58
2.3.4 Método de aceptación-rechazo	59
2.4 Generación a partir de distribuciones de uso común	62
2.4.1 Generación de variables aleatorias continuas	62
2.4.2 Generación de variables aleatorias discretas	67
2.5 Generación de vectores aleatorios	70
2.5.1 Método de aceptación-rechazo de vectores	71
2.5.2 Generación de variables a partir de una distribución multinormal	72
2.5.3 Generación de vectores aleatorios uniformes sobre un símplex	73
2.5.4 Generación de vectores aleatorios distribuidos uniformemente sobre una hipérbola e hiperesfera unitarias	74
2.5.5 Generación de vectores aleatorios uniformemente distribuidos dentro de un hiperelipsoide	75
2.6 Generación de procesos de Poisson	75
2.7 Generación de cadenas de Markov y procesos de salto de Markov	77
2.7.1 Paseo aleatorio en un gráfico	78
2.7.2 Generación de procesos de salto de Markov	79
2.8 Generación de procesos gaussianos	80

2.9 Generación de procesos de difusión	81
2.10 Generación de problemas de permutaciones aleatorias	83 85
Referencias	89
3 Simulación de Sistemas de Eventos Discretos	91
3.1 Introducción	91
3.2 Modelos de simulación	92
3.2.1 Clasificación de los modelos de simulación	94
3.3 Reloj de simulación y lista de eventos para DEDS	95
3.4 Simulación de eventos discretos	97
3.4.1 Cola en tándem	97
3.4.2 Problemas del problema del reparador	101 103
Referencias	106
4 Análisis estadístico de sistemas de eventos discretos	107
4.1 Introducción	107
4.2 Estimadores e Intervalos de Confianza	108
4.3 Modelos de simulación estática	110
4.4 Modelos de simulación dinámica	112
4.4.1 Simulación de horizonte finito	114
4.4.2 Simulación de estado estacionario	114
4.5 Método de arranque	126
Problemas	127
Referencias	130
5 Controlando la varianza	133
5.1 Introducción	133
5.2 Variables aleatorias comunes y antitéticas	134
5.3 Variables de Control	137
5.4 Montecarlo condicional	139
5.4.1 Reducción de varianza para modelos de confiabilidad	141
5.5 Muestreo estratificado	144
5.6 Montecarlo multinivel	146
5.7 Muestreo de importancia	149
5.7.1 Muestras ponderadas	149
5.7.2 Método de minimización de varianza	150
5.7.3 Método de entropía cruzada	154
5.8 Muestreo Secuencial de Importancia	159
5.9 Remuestreo secuencial de importancia	165

5.10 Filtrado no lineal para modelos ocultos de Markov	167
5.11 Método de razón de verosimilitud transformada	171
5.12 Prevención de la degeneración de los problemas de muestreo de importancia	174
Referencias	179
	184
6 Cadena Markov Montecarlo	187
6.1 Introducción	187
6.2 Algoritmo de Metropolis-Hastings	188
6.3 Muestreador Hit-and-Run	193
6.4 Muestreador de Gibbs	194
6.5 Modelos de Ising y Potts	197
6.5.1 Modelo de salida	197
6.5.2 Modelo de macetas	198
6.6 Estadísticas bayesianas	200
6.7 Otros muestreadores de Markov	202
6.7.1 Muestreador de rebanadas	204
6.7.2 Muestreador de salto reversible	205
6.8 Recocido simulado	208
6.9 Muestreo perfecto	212
Problemas	214
Referencias	219
7 Análisis de sensibilidad y optimización de Monte Carlo	221
7.1 Introducción	221
7.2 Método de función de puntuación para el análisis de sensibilidad de DESS	224
7.3 Optimización basada en simulación de DESS	231
7.3.1 Aproximación estocástica	232
7.3.2 Método de la contraparte estocástica	237
7.4 Análisis de sensibilidad de problemas DEDS	246
Referencias	252
	255
8 Método de entropía cruzada	257
8.1 Introducción	257
8.2 Estimación de probabilidades de eventos raros	258
8.2.1 Problema de búsqueda de raíces	267
8.2.2 Método de detección de eventos raros	268
8.2.3 Método CE combinado con muestreo de la distribución de varianza cero	271
8.3 Método CE para Optimización	272

8.4 Problema de corte máximo	276
8.5 Problema de partición	282
8.5.1 Complejidad computacional empírica	283
8.6 Problema del viajante de comercio	283
8.6.1 Gráficos incompletos	288
8.6.2 Colocación de nodos	289
8.6.3 Estudios de casos	290
8.7 Optimización continua	291
8.8 Optimización ruidosa	292
8.9 Método MinxEnt	294
Problemas	298
Referencias	303
9 Método de división	307
9.1 Introducción	307
9.2 Recuento de caminatas autoevasivas mediante división	308
9.3 Partición con un factor de partición fijo	310
9.4 Dividir con esfuerzo fijo	313
9.5 División generalizada	314
9.6 División adaptativa	318
9.7 Aplicación de la división a la confiabilidad de la red	321
9.8 Aplicaciones al conteo	322
9.9 Estudios de casos para contar con división	325
9.9.1 Problema de satisfacción (SAT)	325
9.9.2 Conjuntos independientes	330
9.9.3 Coincidencias Perfectas Permanentes y Contables	332
9.9.4 Tablas de contingencia binarias	334
9.9.5 Coloreado de vértices	336
9.10 División como método de muestreo	337
9.11 División para optimización	340
9.11.1 Problemas de optimización	343
continua	344
Referencias	348
10 Método de enumeración estocástica	351
10.1 Introducción	351
10.2 Búsqueda de árboles y conteo de árboles	352
10.3 Algoritmo de Knuth para estimar el costo de una	355
10.4 enumeración estocástica de árbol	357
10.4.1 Combinación de SE con	359
10.5 Oracle Aplicación de SE al conteo	360

10.5.1	Contando el Número de Caminos en una Red	360
10.5.2	Contando SATs	363
10.5.3	Contar el número de coincidencias perfectas en un gráfico bipartito	366
10.6	Aplicación de SE a la Confiabilidad de la Red	368
10.6.1	Resultados Numéricos	370
	Problemas	373
	Referencias	375
Apéndice		377
A.1	Método de la raíz cuadrada de Cholesky	377
A.2	Muestreo exacto de una distribución de Bernoulli condicional	378
A.3	Familias exponenciales	379
A.4	Análisis de sensibilidad	382
A.4.1	Resultados de convexidad	383
A.4.2	Resultados de monotonidad	384
A.5	Un algoritmo CE simple para optimizar la función de picos Filtro	385
A.6	de Kalman de tiempo discreto	385
A.7	Complejidad del problema de	387
A.8	interrupción de Bernoulli	389
A.8.1	Complejidad de algoritmos de eventos raros	389
A.8.2	Complejidad de algoritmos aleatorios: FPRAS y FPAUS	390
A.8.3	SAT en CNF	394
A.8.4	Complejidad de los problemas de programación estocástica	395
	Problemas	402
	Referencias	403
abreviaciones y acronimos		405
Lista de Símbolos		407
	Índice	409

PREFACIO

Desde la publicación en 2008 de la segunda edición de *Simulación y el Método Monte Carlo*, se han producido cambios significativos en el campo de la simulación Monte Carlo. Esta tercera edición ofrece una descripción completa y totalmente actualizada de los principales temas de la simulación Monte Carlo.

El libro se basa en un curso de pregrado sobre métodos de Monte Carlo impartido en el Instituto de Tecnología de Israel (Technion) y la Universidad de Queensland durante los últimos cinco años. Está dirigido a una amplia audiencia de estudiantes de ingeniería, ciencias físicas y de la vida, estadística, informática, matemáticas y simplemente cualquier persona interesada en utilizar la simulación de Monte Carlo en su estudio o trabajo. Nuestro objetivo es proporcionar una introducción accesible a los métodos modernos de Monte Carlo, centrándonos en los conceptos principales, al mismo tiempo que proporcionamos una base sólida para la resolución de problemas. Por esta razón, la mayoría de las ideas se presentan y explican a través de ejemplos concretos, algoritmos y experimentos.

Aunque suponemos que el lector tiene cierta formación matemática básica, como un curso elemental de probabilidad y estadística, en el capítulo 1 repasamos los conceptos básicos de probabilidad, procesos estocásticos, teoría de la información y optimización convexa.

En una simulación estocástica típica, la aleatoriedad se introduce en los modelos de simulación a través de variables aleatorias independientes distribuidas uniformemente. Estas variables aleatorias se utilizan luego como bloques de construcción para simular sistemas estocásticos más generales. El capítulo 2 trata de la generación de tales números aleatorios, variables aleatorias y procesos estocásticos.

Muchos sistemas complejos del mundo real se pueden modelar como sistemas de eventos discretos. Los ejemplos de sistemas de eventos discretos incluyen sistemas de tráfico, fabricación flexible

sistemas informáticos, sistemas de comunicación informática, sistemas de inventario, líneas de producción, sistemas de vida útil coherente, redes PERT y redes de flujo. El comportamiento de tales sistemas se identifica a través de una secuencia de "eventos" discretos que hacen que el sistema cambie de un "estado" a otro. Discutimos cómo modelar tales sistemas en una computadora en el Capítulo 3.

El Capítulo 4 trata el análisis estadístico de los datos de salida de los modelos de simulación estáticos y dinámicos. La principal diferencia es que los primeros no evolucionan en el tiempo mientras que los segundos sí. Para los modelos dinámicos, distinguimos entre simulaciones de horizonte finito y de estado estacionario. También se analizan dos métodos populares para estimar las medidas de rendimiento en estado estacionario: las medias por lotes y los métodos regenerativos.

El capítulo 5 trata de las técnicas de reducción de la varianza en la simulación de Monte Carlo, como números aleatorios antitéticos y comunes, variables aleatorias de control, Monte Carlo condicional, muestreo estratificado y muestreo de importancia. Usando el muestreo de importancia, a menudo se puede lograr una reducción sustancial (a veces dramática) de la varianza, en particular cuando se estiman probabilidades de eventos raros. Al tratar con el muestreo de importancia, presentamos dos enfoques alternativos, llamados *el minimización de la varianzay el entropía cruzada* métodos. Se presta especial atención a los algoritmos de muestreo de importancia en los que los caminos se generan de forma secuencial. Se obtienen mejoras adicionales de dichos algoritmos mediante el remuestreo de rutas exitosas, lo que da lugar a *remuestreo de importancia secuencial* algoritmos. Ilustramos su uso a través de un ejemplo de filtrado no lineal. Además, este capítulo contiene dos nuevos métodos basados en muestreo de importancia, llamados *el método de razón de verosimilitud transformada* y *el método de selección* para la reducción de la varianza. El primero presenta una forma simple, conveniente y unificadora de construir estimadores de muestreo de importancia eficientes, mientras que el segundo asegura la reducción de la dimensionalidad de la densidad de muestreo de importancia. Esto se logra mediante la identificación (descarte) de los parámetros más importantes (cuello de botella) que se utilizarán en la distribución de muestreo de importancia. Como resultado, la precisión del estimador de muestreo de importancia aumenta sustancialmente.

El capítulo 6 da un tratamiento conciso del genérico *Cadena de Markov Montecarlo* (MCMC) método para *aproximadamente* generar muestras a partir de una distribución arbitraria. Analizamos el algoritmo clásico de Metropolis-Hastings y el muestreador de Gibbs. En el primero, se simula una cadena de Markov tal que su distribución estacionaria coincide con la distribución objetivo, mientras que en el segundo, la cadena de Markov subyacente se construye sobre la base de una secuencia de distribuciones condicionales. También tratamos las aplicaciones de MCMC en las estadísticas bayesianas y explicamos cómo se usa MCMC para tomar muestras de la distribución de Boltzmann para los modelos de Ising y Potts, que se usan ampliamente en la mecánica estadística. Además, mostramos cómo se usa MCMC en el método de recocido simulado para encontrar el mínimo global de una función multiextremal. *rodajay el salto reversible* muestreadores

El Capítulo 7 trata sobre el análisis de sensibilidad y la optimización Monte Carlo de sistemas simulados. Debido a su complejidad, la evaluación del desempeño de los sistemas de eventos discretos generalmente se estudia mediante simulación, y la simulación a menudo se asocia con la estimación de la función de desempeño con respecto a algunos parámetros controlables. El análisis de sensibilidad se ocupa de evaluar las sensibilidades (gradientes, hessianas, etc.) de la función de rendimiento con respecto a los parámetros del sistema.

Esto proporciona una guía para las decisiones operativas y para seleccionar los parámetros del sistema que optimizan las medidas de rendimiento. La optimización de Monte Carlo se ocupa de resolver programas estocásticos, es decir, problemas de optimización donde la función objetivo y algunas de las restricciones son desconocidas y deben obtenerse mediante simulación. Nos ocupamos del análisis de sensibilidad y la optimización de modelos estáticos y dinámicos. Presentamos al célebre *método de puntuación* para el análisis de sensibilidad, y dos métodos alternativos para la optimización de Monte Carlo, el llamado *aproximación estocástica* y *contraparte estocástica* métodos. En particular, en el último método, mostramos cómo usando un solo experimento de simulación uno puede aproximarse con bastante precisión a la verdadera solución óptima desconocida del programa determinista original.

El capítulo 8 trata de la *entropía cruzada* (CE), que fue introducido por el primer autor en 1997 como un algoritmo adaptativo para la estimación de eventos raros utilizando una técnica de minimización de entropía cruzada. Pronto se dio cuenta de que las ideas subyacentes tenían un rango de aplicación mucho más amplio que solo en la simulación de eventos raros: podrían adaptarse fácilmente para abordar problemas de optimización combinatoria y multiextremal bastante generales, incluidos muchos problemas asociados con algoritmos de aprendizaje y computación neuronal. Brindamos una introducción gradual al método CE y mostramos su elegancia y versatilidad. En particular, presentamos un algoritmo CE general para la estimación de probabilidades de eventos raros y luego lo modificamos ligeramente para resolver problemas de optimización combinatoria. Discutimos las aplicaciones del método CE a varios problemas de optimización combinatoria, como el problema del corte máximo y el problema del viajante de comercio, y proporciona resultados numéricos de apoyo sobre su eficacia. Debido a su versatilidad, manejabilidad y simplicidad, el método CE tiene potencialmente una amplia gama de aplicaciones, por ejemplo, en biología computacional, alineación de secuencias de ADN, teoría de grafos y programación. En los últimos 10 años se han escrito cientos de artículos sobre la teoría y las aplicaciones de la CE. Para más detalles ver el sitio En los últimos 10 años se han escrito cientos de artículos sobre la teoría y las aplicaciones de la CE. Para más detalles ver el sitio En los últimos 10 años se han escrito cientos de artículos sobre la teoría y las aplicaciones de la CE. Para más detalles ver el sitio www.cemethod.org, nuestro libro *El método de entropía cruzada: un enfoque unificado para la optimización combinatoria, la simulación Monte-Carlo y el aprendizaje automático* (Springer, 2004), y busque en la wikipedia bajo "método de entropía cruzada". El capítulo concluye con una discusión del programa de optimización de mínima entropía cruzada (MinxEnt).

El capítulo 9 presenta la *terrible* método, que utiliza un plan de muestreo secuencial para descomponer un problema "difícil" en una secuencia de problemas "fáciles". El método se diseñó originalmente para la simulación de eventos raros, pero se ha convertido en un algoritmo de "partículas MCMC" muy versátil que se puede utilizar para la estimación, optimización y muestreo de eventos raros. El capítulo presenta varios algoritmos de división para modelos de simulación dinámicos y estáticos, y demuestra cómo se pueden usar para (1) estimar probabilidades de eventos raros, (2) resolver problemas de conteo difíciles, (3) encontrar soluciones a problemas de optimización desafiantes y (4) muestra de distribuciones de probabilidad complicadas. El capítulo presenta una amplia variedad de estudios de casos y experimentos numéricos, lo que demuestra la eficacia del método.

Muchos problemas combinatorios se pueden formular en términos de buscar o contar el costo total de un árbol. El capítulo 10 presenta un nuevo Monte Carlo llamado *enumeración estocástica* (SE) que es muy adecuado para resolver este tipo de problemas mediante la generación de rutas aleatorias a través del árbol de forma paralela. El algoritmo SE puede verse como un método de muestreo de importancia secuencial en un "hiper-árbol" cuyos vértices son conjuntos de vértices del árbol original. Al combinar SE con algoritmos rápidos de decisión polinomial, mostramos cómo se puede usar para contar problemas #P-completos, como

como el número de asignaciones de satisfacibilidad, el número de rutas en una red general y el número de coincidencias perfectas en un gráfico. La utilidad del método se ilustra a través de una serie de ejemplos numéricos.

El apéndice presenta una variedad de temas complementarios, incluida una breve introducción a las familias exponenciales, el filtro de Kalman de tiempo discreto y el método de la raíz cuadrada de Cholesky. También se analiza la complejidad computacional de los algoritmos aleatorios. Al final de cada capítulo se proporciona una amplia gama de ejercicios.

Además de dos capítulos completamente nuevos (Capítulos 9 y 10), esta tercera edición ofrece actualizaciones sustanciales sobre una variedad de temas. El material sobre la generación de números aleatorios se ha revisado exhaustivamente al incluir generadores recursivos múltiples combinados de última generación y el Mersenne Twister. El material sobre la generación de procesos estocásticos se ha ampliado al incluir la simulación de procesos gaussianos, movimiento browniano y procesos de difusión. El capítulo de reducción de la varianza ahora incluye una discusión del novedoso método Monte Carlo multinivel. Nuestro tratamiento de la importancia secuencial se ha modificado significativamente al enfatizar la importancia del remuestreo de importancia. Esta adición también prepara para el enfoque de partículas MCMC en el nuevo capítulo de división. El capítulo de entropía cruzada se mejora aún más al agregar nuevos conocimientos sobre la degeneración de la razón de verosimilitud, lo que lleva al algoritmo CE mejorado de un solo nivel. Se han añadido veinticinco preguntas más, junto con sus soluciones en *el manual de soluciones en línea* que acompaña a este libro. Finalmente, para facilitar su implementación, la mayoría de los algoritmos han sido (re)escritos en pseudocódigo con control de flujo.

REUVEN RUBINSTEIN Y DFASTIDIARKROSA

Haifa y Brisbane

julio de 2016

EXPRESIONES DE GRATITUD

Agradecemos a todos los que contribuyeron a este libro. Robert Smith y Zelda Zabinski leyeron y proporcionaron sugerencias útiles sobre el Capítulo 6. Alex Shapiro dio una descripción detallada de la complejidad de los problemas de programación estocástica (Sección A.8.4). Pierre L'Ecuyer comunicó el contraejemplo en la Sección 9.10 y Zdravko Botev sugirió el ejemplo multinivel de Monte Carlo en la Sección 5.6. Jim Spall tuvo la amabilidad de proporcionar comentarios sobre la segunda edición, que hemos incorporado en la tercera edición. Josh Chan, Tom Taimre y Zdravko Botev ayudaron a corregir el nuevo material de este libro, lo cual es muy apreciado.

Estamos especialmente agradecidos a los muchos estudiantes de pregrado y posgrado del Technion y de la Universidad de Queensland que ayudaron a hacer posible este libro y cuyas valiosas ideas y experimentos fueron extremadamente alentadores y motivadores. Qibin Duan, Morgan Grant, Robert Salomone, Rohan Shah y Erli Wang leyeron el nuevo material, resolvieron los nuevos ejercicios y brindaron excelentes comentarios. Un agradecimiento especial a Slava Vaisman por su ayuda en muchos problemas computacionales que encontramos, y por las muchas discusiones fructíferas que tuvimos sobre los métodos de enumeración estocástica y división.

Este libro fue apoyado por el Australian Research Council *Centro de Excelencia para las Fronteras Matemáticas y Estadísticas*, con número de concesión CE140100049.

RYR, DPK

CAPÍTULO 1

PRELIMINARES

1.1 INTRODUCCIÓN

El propósito de este capítulo es revisar algunos hechos básicos de la probabilidad, la teoría de la información y la optimización. En particular, las secciones 1.2 a 1.11 resumen los puntos principales de la teoría de la probabilidad. Las secciones 1.12 a 1.14 describen varios procesos estocásticos fundamentales, como los procesos de Poisson, Markov y Gauss. Los elementos de la teoría de la información se dan en la Sección 1.15, y la Sección 1.16 concluye con un resumen de la teoría de la optimización convexa.

1.2 EXPERIMENTOS AL AZAR

La noción básica en la teoría de la probabilidad es la de un *experimento aleatorio*: un experimento cuyo resultado no se puede determinar de antemano. El ejemplo más fundamental es el experimento en el que se lanza una moneda justa varias veces. Por simplicidad, suponga que la moneda se lanza tres veces. El *espacio muestral*, denotado Ω , es el conjunto de todos los resultados posibles del experimento. En este caso Ω tiene ocho resultados posibles:

$$\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\},$$

donde, por ejemplo, HTH significa que el primer lanzamiento es cara, el segundo cruz y el tercero cara.

Los subconjuntos del espacio muestral se llaman *eventos*. Por ejemplo, el evento A que el tercer lanzamiento es cara es

$$A = \{HHH, HTH, THH, TTH\}.$$

Decimos ese evento A ocurre si el resultado del experimento es uno de los elementos en A . Dado que los eventos son conjuntos, podemos aplicarles las operaciones habituales de conjuntos. Por ejemplo, el evento $A \cup B$, llamó al *Unión de A y B* , es el evento que A o B o ambos ocurren, y el evento $A \cap B$, llamó al *intersección de A y B* , es el evento que A y B ambos ocurren. Una notación similar es válida para uniones e intersecciones de más de dos eventos. El evento A^c , llamó al *complementario de A* , es el evento que A no se produce. dos eventos A y B que no tienen resultados en común, es decir, su intersección está vacía, se llaman *desarticulados* eventos. El paso principal es especificar la probabilidad de cada evento.

Definición 1.2.1 (Probabilidad) Una *probabilidad* P es una regla que asigna un número $0 \leq P(A) \leq 1$ para cada evento A , tal que $P(\Omega) = 1$, y tal que para cualquier sucesión A_1, A_2, \dots de eventos disjuntos

$$P\left(\bigcup_i A_i\right) = \sum_i P(A_i). \quad (1.1)$$

La ecuación (1.1) se conoce como la *regla de la suma* de probabilidad. Afirma que si un evento puede ocurrir de varias maneras diferentes, pero no simultáneamente, la probabilidad de ese evento es simplemente la suma de las probabilidades de los eventos que lo componen.

Para el experimento del lanzamiento justo de una moneda, la probabilidad de cualquier evento se da fácilmente. Es decir, debido a que la moneda es justa, cada uno de los ocho resultados posibles es igualmente probable, de modo que $P(\{HHH\}) = \dots = P(\{TTT\}) = 1/8$. Desde cualquier evento A es la unión de los eventos "elementales" $\{HHH\}, \dots, \{TTT\}$, la regla de la suma implica que

$$P(A) = \frac{|A|}{|\Omega|} \quad (1.2)$$

dónde $|A|$ denota el número de resultados en A y $|\Omega| = 8$. En términos más generales, si un experimento aleatorio tiene un número finito de resultados igualmente probables, la probabilidad es siempre de la forma (1.2). En ese caso, el cálculo de probabilidades se reduce a contar.

1.3 PROBABILIDAD CONDICIONAL E INDEPENDENCIA

¿Cómo cambian las probabilidades cuando sabemos que algún evento $B \subset \Omega$ ha ocurrido? Dado que el resultado está en B , el evento A ocurrirá si y sólo si $A \cap B$ ocurre, y la probabilidad relativa de A ocurriendo es por lo tanto $P(A \cap B)/P(B)$. Esto conduce a la definición de la *probabilidad condicional* de A dado B :

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. \quad (1.3)$$

Por ejemplo, supongamos que lanzamos una moneda justa tres veces. Dejar B sea el evento de que el número total de caras sea dos. La probabilidad condicional del evento A que el primer lanzamiento es cara, dado que B ocurre, es $(2/8)/(3/8) = 2/3$.

Reescribiendo (1.3) e intercambiando el papel de A y B da la relación $\text{PAGS}(A \cap B) = \text{PAGS}(A) \text{PAGS}(B/A)$. Esto se puede generalizar fácilmente a la *regla del productode* probabilidad, que establece que para cualquier secuencia de eventos $A_1, A_2, \dots, A_{norte}$,

$$\text{PAGS}(A_1 \cdots A_{norte}) = \text{PAGS}(A_1) \text{PAGS}(A_2/A_1) \text{PAGS}(A_3/A_1 A_2) \cdots \text{PAGS}(A_{norte}/A_1 \cdots A_{n-1}), \quad (1.4)$$

usando la abreviatura $A_1 A_2 \cdots A_k \equiv A_1 \cap A_2 \cap \cdots \cap A_k$.

Suponer que $B_1, B_2, \dots, B_{norte}$ es un *dividirde* Ω . Eso es, $B_1, B_2, \dots, B_{norte}$ son disjuntos y su unión es Ω . Entonces, por la regla de la suma, $\text{PAGS}(A) = \sum_{i=1}^{norte} \text{PAGS}(A \cap B_i)$ y por lo tanto, por la definición de probabilidad condicional, tenemos la *ley de probabilidad total*:

$$\text{PAGS}(A) = \sum_{i=1}^{norte} \text{PAGS}(A/B_i) \text{PAGS}(B_i). \quad (1.5)$$

Combinando esto con la definición de probabilidad condicional da *regla de bayes*:

$$\text{PAGS}(B_i/A) = \frac{\text{PAGS}(A/B_i) \text{PAGS}(B_i)}{\sum_{j=1}^{norte} \text{PAGS}(A/B_j) \text{PAGS}(B_j)}. \quad (1.6)$$

La independencia es de crucial importancia en probabilidad y estadística. Flojamente hablando, modela la falta de información entre eventos. dos eventos A y B se dice que son *independientes* si el conocimiento que B ha ocurrido no cambia la probabilidad de que A ocurre. Eso es, A y B independiente $\Leftrightarrow \text{PAGS}(A/B) = \text{PAGS}(A)$. Ya que $\text{PAGS}(A/B) = \text{PAGS}(A \cap B)/\text{PAGS}(B)$, una definición alternativa de independencia es

$$A, B \text{ independiente} \Leftrightarrow \text{PAGS}(A \cap B) = \text{PAGS}(A) \text{PAGS}(B).$$

Esta definición cubre el caso en que $B = \emptyset$ (conjunto vacío). Podemos extender esta definición a arbitrariamente muchos eventos.

Definición 1.3.1 (Independencia) Los eventos A_1, A_2, \dots , se dice que son *independientes* si por alguna k y cualquier elección de índices distintos i_1, \dots, i_k ,

$$\text{PAGS}(A_{i_1} \cap A_{i_2} \cap \cdots \cap A_{i_k}) = \text{PAGS}(A_{i_1}) \text{PAGS}(A_{i_2}) \cdots \text{PAGS}(A_{i_k})$$

Observación 1.3.1 En la mayoría de los casos, la independencia de los eventos es una suposición del modelo. Es decir, suponemos que existe un PAGS tal que ciertos eventos son independientes.

■ EJEMPLO 1.1

Lanzamos una moneda sesgada $norte$ veces. Dejar pag ser la probabilidad de cara (para una moneda justa $pag = 1/2$). Dejar A denota el evento de que el i -th tirada da cara, $i = 1, \dots, norte$. Después PAGS debe ser tal que los eventos A_1, \dots, A_{norte} son independientes y $\text{PAGS}(A_i) = pag$ para todos i . Estas dos reglas especifican completamente PAGS . Por ejemplo, la probabilidad de que la primera k los tiros son cara y el último $n-k$ son colas es

$$\begin{aligned} \text{PAGS}(A_1 \cdots A_k A_{k+1} \cdots A_{norte}) &= \text{PAGS}(A_1) \cdots \text{PAGS}(A_k) \text{PAGS}(A_{k+1}) \cdots \text{PAGS}(A_{norte}) \\ &= pag^k (1 - pag)^{n-k}. \end{aligned}$$

1.4 VARIABLES ALEATORIAS Y DISTRIBUCIONES DE PROBABILIDAD

Especificar un modelo para un experimento aleatorio a través de una descripción completa de Ω y PAGES no siempre es conveniente o necesario. En la práctica, sólo estamos interesados en ciertas observaciones (es decir, medidas numéricas) en el experimento. Los incorporamos a nuestro proceso de modelado a través de la introducción de *variables aleatorias*, generalmente indicado con letras mayúsculas de la última parte del alfabeto (p. ej., X, X_1, X_2, \dots, Y, Z).

■ EJEMPLO 1.2

Lanzamos una moneda sesgada *norte* veces, con *pags* la probabilidad de cara. Supongamos que solo estamos interesados en el número de cabezas, digamos X . Tenga en cuenta que X puede tomar cualquiera de los valores de $\{0, 1, \dots, n\}$. la *Distribución de probabilidad de X* es dado por el *fórmula binomial*

$$PAGS(X=k) = \binom{n}{k} pags^k (1-pags)^{n-k}, \quad k=0, 1, \dots, n \quad (1.7)$$

Es decir, por el Ejemplo 1.1, cada evento elemental $\{HTH \dots T\}$ con exactamente k cabezas y $n-k$ cruz tiene probabilidad $pags^k (1-pags)^{n-k}$, y hay $\binom{n}{k}$ tal

La distribución de probabilidad de una variable aleatoria general X —identificando tales probabilidades como $PAGS(X=x)$, $PAGS(a < X < b)$, y así sucesivamente, está completamente especificado por la *función de distribución acumulativa* (*cdf*), definido por

$$F(x) = PAGS(X \leq x), \quad x \in \mathbb{R}.$$

Una variable aleatoria X se dice que tiene una *discreta* distribución si, for algo finito o conjunto contable de valores x_1, x_2, \dots , $PAGS(X=x_i) > 0, i=1, 2, \dots$ y $PAGS(X=x) = 0$ para x no en el conjunto. La función $F(x) = PAGS(X \leq x)$ se llama la *función de probabilidad* (*pmf*) de X —pero véase la observación 1.4.1.

■ EJEMPLO 1.3

Lance dos dados justos y deje *METRO* ser el mayor valor nominal que muestra. El pmf de *METRO* es dado por

<i>metro</i>	1	2	3	4	5	6	Σ
$F(\text{metro})$	$\frac{1}{36}$	$\frac{3}{36}$	$\frac{5}{36}$	$\frac{7}{36}$	$\frac{9}{36}$	$\frac{11}{36}$	1

Por ejemplo, para obtener *METRO*=3, ya sea (1,3),(2,3),(3,3),(3,2 o 3,1) tiene que ser lanzado, cada uno de los cuales sucede con probabilidad $1/36$.

Una variable aleatoria X se dice que tiene una *continua* distribución si existe una función positiva f con integral total 1, tal que para todo a, b ,

$$PAGS(a < X < b) = \int_a^b f(t) dt \quad (1.8)$$

La función f se llama *función de densidad de probabilidad* (pdf) de X . Tenga en cuenta que en el caso continuo el cdf está dado por

$$F(x) = \text{PAGS}(X \leq x) = \int_{-\infty}^x f(t) dt,$$

y f es la derivada de F . podemos interpretar $f(x)$ como la “densidad” de probabilidad en x en el sentido de que

$$\text{PAGS}(x \leq X \leq x+h) = \int_x^{x+h} f(t) dt \approx h \cdot f(x).$$

Observación 1.4.1 (Densidad de probabilidad) Tenga en cuenta que hemos utilizado deliberadamente el *mismo* símbolo, f , tanto para pmf como para pdf. Esto se debe a que la pmf y la pdf juegan papeles muy similares y pueden, en una teoría de la probabilidad más avanzada, ser vistas como instancias particulares de la noción general de *densidad de probabilidad*. Para enfatizar este punto de vista, llamaremos f en *ambas cosas* el caso discreto y continuo la pdf o (probabilidad) densidad (función).

1.5 ALGUNAS DISTRIBUCIONES IMPORTANTES

Las tablas 1.1 y 1.2 enumeran una serie de importantes distribuciones continuas y discretas. Usaremos la notación $X \sim F$, $X \sim \text{Dist}$ para significar que X tiene un pdf f , un CDF F o una distribución Dist . A veces escribimos f_X en vez de f para enfatizar que el pdf se refiere a la variable aleatoria X . Note que en la Tabla 1.1, Γ es la función gamma: $\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$, $\alpha > 0$.

Tabla 1.1: Distribuciones continuas comúnmente utilizadas.

Nombre	Notación	$f(x)$	$X \in$	Parámetros
Uniforme	$U[a, b]$	$\frac{1}{b-a}$	$[a, b]$	$a < b$
Normal	$N(\mu, \sigma^2)$	$\frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$	\mathbb{R}	$\sigma > 0, \mu \in \mathbb{R}$
Gama	$\text{Gamma}(\alpha, \lambda)$	$\frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}$	\mathbb{R}_+	$\alpha, \lambda > 0$
Exponencial	$\text{Exp}(\lambda)$	$\lambda e^{-\lambda x}$	\mathbb{R}_+	$\lambda > 0$
Beta	$\text{Beta}(\alpha, \beta)$	$\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$	$[0, 1]$	$\alpha, \beta > 0$
Weibull	$\text{Weib}(\alpha, \lambda)$	$\alpha \lambda^\alpha x^{\alpha-1} e^{-(\lambda x)^\alpha}$	\mathbb{R}_+	$\alpha, \lambda > 0$
Pareto	$\text{Pareto}(\alpha, \lambda)$	$\alpha \lambda^\alpha x^{-(\alpha+1)}$	\mathbb{R}_+	$\alpha, \lambda > 0$

Tabla 1.2: Distribuciones discretas comúnmente utilizadas.

Nombre	Notación	$f(x)$	$X \in$	Parámetros
Bernoulli	$Ber(pags)$	$\binom{pags}{x} pags^x (1-pags)^{1-x}$	$\{0, 1\}$	$0 < pags < 1$
Binomio	$Comp(norte, pags)$	$\binom{norte}{x} pags^x (1-pags)^{norte-x}$	$\{0, 1, \dots, norte\}$	$0 < pags < 1, norte \in \mathbb{N}$
Discreto uniforme	$DU\{1, \dots, n\}$	$\frac{1}{norte}$	$\{1, \dots, n\}$	$norte \in \{1, 2, \dots\}$
Geométrico	$GRAMO(pags)$	$pags(1-pags)^{x-1}$	$\{1, 2, \dots\}$	$0 < pags < 1$
veneno	$Poi(\lambda)$	$\frac{e^{-\lambda} \lambda^x}{x!}$	\mathbb{N}	$\lambda > 0$

1.6 EXPECTATIVA

A menudo es útil considerar diferentes tipos de características numéricas de una variable aleatoria. Una de esas cantidades es la expectativa, que mide el valor medio de la distribución.

Definición 1.6.1 (Expectativa) Dejar X ser una variable aleatoria con pdf f . Los *expectativa* (o valor esperado o media) de X , denotado por $MI[X]$ (o algunas veces m), es definido por

$$MI[X] = \begin{cases} \sum x f(x) & \text{caso discreto,} \\ \int_{-\infty}^{\infty} x f(x) dx & \text{caso continuo.} \end{cases}$$

Si X es una variable aleatoria, entonces una función de X , como $h(X)$, es de nuevo una variable aleatoria. Además, el valor esperado de una función de X es simplemente un promedio ponderado de los posibles valores que puede tomar esta función. Es decir, para cualquier función real h

$$MI[h(X)] = \begin{cases} \sum h(x) f(x) & \text{caso discreto,} \\ \int_{-\infty}^{\infty} h(x) f(x) dx & \text{caso continuo.} \end{cases}$$

Otra cantidad útil es la varianza, que mide la extensión o dispersión de la distribución.

Definición 1.6.2 (Varianza) Los *diferenciade* una variable aleatoria X , denotada por $Var(X)$ (o algunas veces σ^2), es definido por

$$Var(X) = MI[(X - MI[X])^2] = MI[X^2] - (MI[X])^2.$$

La raíz cuadrada de la varianza se llama *Desviación Estándar*. La tabla 1.3 enumera las expectativas y las varianzas de algunas distribuciones conocidas.

Tabla 1.3: Expectativas y varianzas para algunas distribuciones conocidas.

Dist.	MI[X]	Var(X)	Dist.	MI[X]	Var(X)
Compartimiento(norte, pag) notario público notario público(1-pags)			Gama(α, λ)	$\frac{\alpha}{\lambda}$	$\frac{\alpha}{\lambda^2}$
GRAMQ(pags)	$\frac{1}{pags}$	$\frac{1-pags}{pags}$	NORTE(μ, σ^2)	m	σ^2
Poi(λ)	λ	λ	Beta(α, β)	$\frac{\alpha}{\alpha+\beta}$	$\frac{\alpha\beta}{(\alpha+\beta)^2(1+\alpha+\beta)}$
tu(α, β)	$\frac{\alpha+\beta}{2}$	$\frac{(\beta-\alpha)^2}{12}$	Weib(α, λ)	$\frac{\Gamma(1/\alpha)}{\alpha\lambda}$	$\frac{2\Gamma(2/\alpha)}{\alpha} - \left(\frac{\Gamma(1/\alpha)}{\alpha\lambda}\right)^2$
Exp(λ)	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$			

La media y la varianza no dan, en general, suficiente información para especificar completamente la distribución de una variable aleatoria. Sin embargo, pueden proporcionar límites útiles. Discutimos dos de tales límites. Suponer X solo puede tomar valores no negativos y tiene pdf f . Para cualquier $X > 0$, podemos escribir

$$MI[X] = \int_0^\infty t \cdot f(t) dt + \int_X^\infty t \cdot f(t) dt - \int_X^\infty t \cdot f(t) dt$$

de donde se sigue la *Desigualdad de Markov*: si $X > 0$, entonces para todos $X > 0$,

$$PAGS(X-X) - \frac{MI[X]}{X}. \quad (1.9)$$

Si también conocemos la varianza de una variable aleatoria, podemos dar un límite más estrecho. Es decir, para cualquier variable aleatoria X con media m y varianza σ^2 , tenemos

$$PAGS(|X - \mu| - X) - \frac{\sigma^2}{X^2}. \quad (1.10)$$

Esto se llama el *Desigualdad de Chebyshev*. La demostración es la siguiente: Sea $D_2 = (X - \mu)^2$; luego, por la desigualdad de Markov (1.9) y la definición de la varianza,

$$PAGS(D_2 - X_2) - \frac{\sigma^2}{X_2}.$$

Además, tenga en cuenta que el evento $\{D_2 \geq X_2\}$ es equivalente al evento $\{|X - \mu| \geq \sqrt{X_2}\}$, de modo que se sigue (1.10).

1.7 DISTRIBUCIONES CONJUNTAS

A menudo, un experimento aleatorio se describe mediante más de una variable aleatoria. La teoría para múltiples variables aleatorias es similar a la de una sola variable aleatoria.

Dejar $X_1, \dots, X_{\text{norte}}$ Ser variables aleatorias que describen algún experimento aleatorio. Podemos acumularlos en un *vector aleatorio* $\mathbf{X} = (X_1, \dots, X_{\text{norte}})$. Más generalmente, una colección $\{X_t, t \in T\}$ de variables aleatorias se llama *Proceso estocástico*. El conjunto T se llama el *conjunto de parámetros* o *conjunto de índices* del proceso. Puede ser discreto (por ejemplo, $\text{norte} = \{1, \dots, 10\}$) o continuo (p. ej., $\mathbb{R}_+ = [0, \infty)$ o $[1, 10]$). El conjunto de valores posibles para el proceso estocástico se denomina *espacio de Estados*.

La distribución conjunta de $X_1, \dots, X_{\text{norte}}$ es especificado por el *cdf conjunto*

$$F_{\mathbf{X}}(X_1, \dots, X_{\text{norte}}) = \text{PAGS}(X_1 - X_1, \dots, X_{\text{norte}} - X_{\text{norte}}).$$

los *pdf conjunto* F viene dada, en el caso discreto, por $F_{\mathbf{X}}(X_1, \dots, X_{\text{norte}}) = \text{PAGS}(X_1 = X_1, \dots, X_{\text{norte}} = X_{\text{norte}})$, y en el caso continuo F es tal que

$$\text{PAGS}(\mathbf{X} \in B) = \int_B F_{\mathbf{X}}(X_1, \dots, X_{\text{norte}}) dX_1 \dots dX_{\text{norte}}$$

para cualquier región (medible) B en $\mathbb{R}_{\text{norte}}$. Los pdf marginales se pueden recuperar del pdf conjunto por integración o suma. Por ejemplo, en el caso de un vector aleatorio continuo (X, Y) con pdf conjunto F , el pdf F_X de X se encuentra como

$$F_X(X) = \int F_X(X, y) dy.$$

Suponer que X y Y son ambos discretos o ambos continuos, con pdf conjunto F , y supongamos que $F_X(X) > 0$. Entonces el *pdf condicional* de Y dado $X = X$ es dado por

$$F_{Y|X}(y | x) = \frac{F_X(X, y)}{F_X(X)} \quad \text{para todos } y.$$

El correspondiente *expectativa condicional* (en el caso continuo)

$$\text{MI}[Y | X = X] = \int y f_{Y|X}(y | x) dy.$$

Tenga en cuenta que $\text{MI}[Y | X = X]$ es una función de X , decir $h(X)$. La variable aleatoria correspondiente $h(X)$ se escribe como $\text{MI}[Y | X]$. Se puede demostrar (ver, por ejemplo, [3]) que su expectativa es simplemente la expectativa de Y , eso es,

$$E[E[Y | X]] = \text{MI}[Y]. \quad (1.11)$$

Cuando la distribución condicional de Y dado X es idéntico al de Y , X y Y se dice que son independientes. Más precisamente:

Definición 1.7.1 (Variables aleatorias independientes) Las variables aleatorias $X_1, \dots, X_{\text{norte}}$ son llamados *independientes* si para todos los eventos $\{X_i \in A_i\}$ con $A_i \subset \mathbb{R}$, $i = 1, \dots, \text{norte}$,

$$\text{PAGS}(X_1 \in A_1, \dots, X_{\text{norte}} \in A_{\text{norte}}) = \text{PAGS}(X_1 \in A_1) \cdots \text{PAGS}(X_{\text{norte}} \in A_{\text{norte}}).$$

Una consecuencia directa de la definición anterior de independencia es que las variables aleatorias $X_1, \dots, X_{\text{norte}}$ con pdf conjunto F (discretas o continuas) son independientes si y sólo si

$$F(X_1, \dots, X_{\text{norte}}) = F_X(X_1) \cdots F_X(X_{\text{norte}}) \quad (1.12)$$

para todos $X_1, \dots, X_{\text{norte}}$, donde $\{F_X\}$ son los pdf marginales.

■ EJEMPLO 1.4 Secuencia de Bernoulli

Considere el experimento en el que lanzamos una moneda sesgada n veces, con probabilidad p de cabezas. Podemos modelar este experimento de la siguiente manera. Para $i = 1, \dots, n$, dejar X_i ser el resultado de la i -th lanzamiento: $\{X_i = 1\}$ significa cara (o éxito), $\{X_i = 0\}$ significa cruz (o fracaso). También, deja

$$P(X_i = 1) = p, \quad P(X_i = 0) = 1 - p, \quad i = 1, 2, \dots, n$$

Por último, suponga que X_1, \dots, X_n son independientes. La secuencia $\{X_i, i = 1, 2, \dots\}$ se llama *unsecuencia de Bernoulli* o *proceso de Bernoulli* con probabilidad de éxito p . Dejar $X = X_1 + \dots + X_n$ sea el número total de éxitos en n lanzamientos.

Denotamos por B el conjunto (fa) de los vectores binarios $\mathbf{X} = (X_1, \dots, X_n)$ tal que $\sum_{i=1}^n X_i = k$. Tenga en cuenta que B posee $\binom{n}{k}$ elementos. ahora tenemos

$$\begin{aligned} P(X = k) &= \sum_{\mathbf{X} \in B} P(X_1 = X_1, \dots, X_n = X_n) \\ &= \sum_{\mathbf{X} \in B} P(X_1 = X_1) \cdots P(X_n = X_n) = \sum_{\mathbf{X} \in B} p^k (1-p)^{n-k} \\ &= \binom{n}{k} p^k (1-p)^{n-k}. \end{aligned}$$

En otras palabras, $X \sim \text{Binomial}(n, p)$. Compare esto con el Ejemplo 1.2.

Observación 1.7.1 Una *infinita* secuencia X_1, X_2, \dots de variables aleatorias se llama *independiente* si para cualquier elección finita de parámetros i_1, i_2, \dots, i_n (ninguno de ellos es igual a 0) las variables aleatorias X_{i_1}, \dots, X_{i_n} son independientes. Muchos modelos probabilísticos involucran variables aleatorias X_1, X_2, \dots que son *independientes e idénticamente distribuidas*, abreviado como *i.i.d.* Usaremos esta abreviatura a lo largo de este libro.

Similar al caso unidimensional, el valor esperado de cualquier función de valor real h de X_1, \dots, X_n es un promedio ponderado de todos los valores que puede tomar esta función. Específicamente, en el caso continuo,

$$E[h(X_1, \dots, X_n)] = \int \cdots \int h(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \cdots dx_n.$$

Como consecuencia directa de las definiciones de expectativa e independencia, tenemos

$$E[a + b_1 X_1 + b_2 X_2 + \cdots + b_n X_n] = a + b_1 m_1 + \cdots + b_n m_n \quad (1.13)$$

para cualquier secuencia de variables aleatorias X_1, X_2, \dots, X_n con expectativas m_1, m_2, \dots, m_n , donde a, b_1, b_2, \dots, b_n son constantes. Del mismo modo, para *independiente* variables aleatorias, tenemos

$$E[X_1 X_2 \cdots X_n] = m_1 m_2 \cdots m_n.$$

El *covarianza* de dos variables aleatorias X y Y con expectativas $E[X] = \mu_X$ y $E[Y] = \mu_Y$, respectivamente, se define como

$$\text{cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)].$$

Esta es una medida de la cantidad de dependencia lineal entre las variables. Una versión escalada de la covarianza viene dada por la *coeficiente de correlación*,

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y},$$

dónde $\sigma_X = \sqrt{\text{Var}(X)}$ y $\sigma_Y = \sqrt{\text{Var}(Y)}$. Se puede demostrar que la correlación siempre se encuentra entre -1 y 1 ; vea el Problema 1.13.

Para facilitar la referencia, la tabla 1.4 enumera algunas propiedades importantes de la varianza y la covarianza. Las pruebas se derivan directamente de las definiciones de covarianza y varianza y las propiedades de la expectativa.

Tabla 1.4: Propiedades de varianza y covarianza.

1	$\text{Var}(X) = \text{MI}[X^2] - (\text{MI}[X])^2$
2	$\text{Var}(aX + b) = a^2 \text{Var}(X)$
3	$\text{cov}(X, Y) = \text{MI}[XY] - \text{MI}[X]\text{MI}[Y]$
4	$\text{cov}(X, Y) = \text{cov}(Y, X)$
5	$\text{cov}(aX + bY, Z) = a \text{cov}(X, Z) + b \text{cov}(Y, Z)$
6	$\text{cov}(X, X) = \text{Var}(X)$
7	$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{cov}(X, Y)$
8	$X \text{ y } Y \text{ independiente} \Rightarrow \text{cov}(X, Y) = 0$

Como consecuencia de las propiedades 2 y 7, para cualquier secuencia de *independiente* aleatorio Variables $X_1, \dots, X_{\text{norte}}$ con variaciones $\sigma_1^2, \dots, \sigma_{\text{norte}}^2$,

$$\text{Var}(a + b_1 X_1 + b_2 X_2 + \dots + b_{\text{norte}} X_{\text{norte}}) = \sigma_1^2 b_1^2 + \dots + \sigma_{\text{norte}}^2 b_{\text{norte}}^2 \tag{1.14}$$

para cualquier elección de constantes $a, b_1, \dots, b_{\text{norte}}$.

Para vectores aleatorios, como $\mathbf{X} = (X_1, \dots, X_{\text{norte}})$, es conveniente escribir las expectativas y covarianzas en notación vectorial.

Definición 1.7.2 (Vector de expectativa y matriz de covarianza) Para cualquier vector aleatorio \mathbf{X} , definimos el *expectativa de vectores* como vector de expectativas

$$\mathbf{m} = (m_1, \dots, m_{\text{norte}}) = (\text{MI}[X_1], \dots, \text{MI}[X_{\text{norte}}]).$$

los *Matriz de covarianza* Σ se define como la matriz cuya (i, j) -ésimo elemento es

$$\text{cov}(X_i, X_j) = \text{MI}[(X_i - m_i)(X_j - m_j)].$$

Si definimos la expectativa de un vector (matriz) como el vector (matriz) de las expectativas, entonces podemos escribir

$$\mathbf{m} = \text{MI}[\mathbf{X}]$$

y

$$\Sigma = \text{MI}[(X - m)(X - m)'].$$

Tenga en cuenta que m y Σ asumen el mismo papel que m y σ^2 en el caso unidimensional.

Observación 1.7.2 Tenga en cuenta que cualquier matriz de covarianza Σ es *simétrica*. De hecho (vea el Problema 1.16), *esemidefinido positivo*, es decir, para cualquier vector (columna) u ,

$$u' \Sigma u \geq 0.$$

1.8 FUNCIONES DE VARIABLES ALEATORIAS

Suponer que X_1, \dots, X_n son medidas de un experimento aleatorio. A menudo sólo estamos interesados en ciertas *funciones* de las medidas en lugar de las medidas individuales. Aquí hay unos ejemplos.

■ EJEMPLO 1.5

Dejar X ser una variable aleatoria continua con pdf F_X . Dejar $Z = aX + b$, donde $a > 0$. Deseamos determinar el pdf F_Z de Z . Suponer que $a > 0$. Tenemos para cualquier z

$$F_Z(z) = \text{PAGS}(Z \leq z) = \text{PAGS}(X \leq (z - b)/a) = F_X((z - b)/a).$$

Diferenciando esto con respecto a z (da $F_Z'(z) = F_X'((z - b)/a) \cdot (1/a)$). Para $a < 0$ obtenemos de manera similar $F_Z(z) = F_X((z - b)/a)$. Así, en general,

$$F_Z'(z) = \frac{1}{|a|} F_X'\left(\frac{z - b}{a}\right). \quad (1.15)$$

■ EJEMPLO 1.6

Generalizando el ejemplo anterior, supongamos que $Z = g(X)$ para alguna función monótonamente creciente g . Para encontrar el pdf de Z de eso de X nosotros primero escribiremos

$$F_Z(z) = \text{PAGS}(Z \leq z) = \text{PAGS}(X \leq g^{-1}(z)) = F_X(g^{-1}(z)),$$

donde g^{-1} es el inverso de g . Diferenciar con respecto a z ahora da

$$F_Z'(z) = F_X'(g^{-1}(z)) \cdot \frac{d}{dz} g^{-1}(z) = \frac{F_X'(g^{-1}(z))}{g'(g^{-1}(z))}. \quad (1.16)$$

Para funciones monótonamente decrecientes, $\frac{d}{dz} g^{-1}(z)$ en la primera ecuación necesita ser reemplazado por su valor negativo.

■ EJEMPLO 1.7 Estadísticas de pedidos

Dejar X_1, \dots, X_n ser una secuencia iid de variables aleatorias con pdf común F_X CDF F . En muchas aplicaciones uno está interesado en la distribución de la

estadísticas de pedidos $X_{(1)}, X_{(2)}, \dots, X_{(norte)}$, donde $X_{(1)}$ es el más pequeño de los $\{X_i, i = 1, \dots, n\}$, $X_{(2)}$ es el segundo más pequeño, y así sucesivamente. el CDF de $X_{(norte)}$ sigue desde

$$PAGS(X_{(norte)} \leq x) = PAGS(X_1 \leq x, \dots, X_{norte} \leq x) = \prod_{i=1}^{norte} PAGS(X_i \leq x) = (F(x))^{norte}.$$

Similarmemente,

$$PAGS(X_{(1)} > x) = PAGS(X_1 > x, \dots, X_{norte} > x) = \prod_{i=1}^{norte} PAGS(X_i > x) = (1 - F(x))^{norte}.$$

Además, debido a que todos los pedidos de X_1, \dots, X_{norte} son igualmente probables, se deduce que el pdf conjunto de la muestra ordenada es, en la cuña $\{(X_1, \dots, X_{norte}) : X_1 \leq X_2 \leq \dots \leq X_{norte}\}$, simplemente $norte!$ veces la densidad conjunta de la muestra desordenada y cero en el resto.

1.8.1 Transformaciones lineales

Dejar $\mathbf{X} = (X_1, \dots, X_{norte})$ sea un vector columna en \mathbb{R}^{norte} y \mathbf{A} una $metro \times norte$ matriz. el mapeo $\mathbf{X} \rightarrow \mathbf{Z}$, con $\mathbf{Z} = \mathbf{A}\mathbf{X}$, se llama una *transformación lineal*. Ahora considere un aleatorio vector $\mathbf{X} = (X_1, \dots, X_{norte})$, y deja

$$\mathbf{Z} = \mathbf{A}\mathbf{X}.$$

Después \mathbf{Z} es un vector aleatorio en \mathbb{R}^{metro} . En principio, si conocemos la distribución conjunta de \mathbf{X} , entonces podemos derivar la distribución conjunta de \mathbf{Z} . Veamos primero cómo se transforman el vector de expectativas y la matriz de covarianza.

Teorema 1.8.1 Si \mathbf{X} tiene un vector de expectativa \mathbf{m}_x y matriz de covarianza Σ_x , entonces el vector de expectativa y la matriz de covarianza de $\mathbf{Z} = \mathbf{A}\mathbf{X}$ son dados por

$$\mathbf{m}_z = \mathbf{A}\mathbf{m}_x \quad (1.17)$$

y

$$\Sigma_z = \mathbf{A}\Sigma_x\mathbf{A}^T. \quad (1.18)$$

Prueba: Tenemos $\mathbf{m}_z = \mathbb{E}[\mathbf{Z}] = \mathbb{E}[\mathbf{A}\mathbf{X}] = \mathbf{A}\mathbb{E}[\mathbf{X}] = \mathbf{A}\mathbf{m}_x$ y

$$\begin{aligned} \Sigma_z &= \mathbb{E}[(\mathbf{Z} - \mathbf{m}_z)(\mathbf{Z} - \mathbf{m}_z)^T] = \mathbb{E}[(\mathbf{A}\mathbf{X} - \mathbf{A}\mathbf{m}_x)(\mathbf{A}\mathbf{X} - \mathbf{A}\mathbf{m}_x)^T] \\ &= \mathbf{A} \mathbb{E}[(\mathbf{X} - \mathbf{m}_x)(\mathbf{X} - \mathbf{m}_x)^T] \mathbf{A}^T = \mathbf{A}\Sigma_x\mathbf{A}^T. \end{aligned}$$

Suponer que \mathbf{A} es una invertible $norte \times norte$ matriz. Si \mathbf{X} tiene una densidad conjunta f_x , ¿cuál es la densidad conjunta f_z de \mathbf{Z} ? Considere la figura 1.1. Para cualquier fijo \mathbf{x} , dejar $\mathbf{z} = \mathbf{A}\mathbf{x}$. Por eso, $\mathbf{x} = \mathbf{A}^{-1}\mathbf{z}$. Considera el $norte$ -cubo dimensional $C = [z_1, z_1 + h] \times \dots \times [z_{norte}, z_{norte} + h]$. Dejar D ser la imagen de C por debajo de \mathbf{A}^{-1} , es decir, el paralelepípedo de todos los puntos \mathbf{x} tal que $\mathbf{A}\mathbf{x} \in C$. Después,

$$PAGS(\mathbf{Z} \in C) = h^{norte} f_z(\mathbf{z}).$$

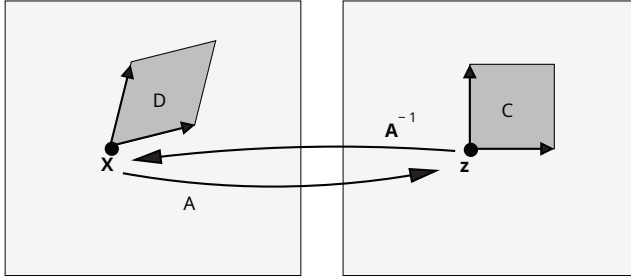


Figura 1.1: Transformación lineal.

Ahora recuerde del álgebra lineal (p. ej., [5]) que cualquier matriz B transforma linealmente un n -rectángulo dimensional con volumen V en un n -paralelepípedo dimensional con volumen $V / |B|$, donde $|B| = |\det(B)|$. De este modo,

$$\text{PAGS}(\mathbf{Z} \in C) = \text{PAGS}(\mathbf{X} \in D) \approx h_{\text{norte}} / |U|^{n-1} / f_{\mathbf{X}}(\mathbf{X}) = h_{\text{norte}} / |A|^{-1} f_{\mathbf{X}}(\mathbf{X}).$$

Alquilar h vamos a 0, obtenemos

$$f_{\mathbf{Z}}(\mathbf{z}) = \frac{f_{\mathbf{X}}(A^{-1}\mathbf{z})}{|A|}, \quad \mathbf{z} \in \mathbb{R}_{\text{norte}}. \quad (1.19)$$

1.8.2 Transformaciones generales

Podemos aplicar un razonamiento similar al anterior para tratar con transformaciones generales $\mathbf{X} \rightarrow \mathbf{g}(\mathbf{X})$, escrito como

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \end{pmatrix} \xrightarrow{\mathbf{g}} \begin{pmatrix} g_1(\mathbf{X}) \\ g_2(\mathbf{X}) \\ \vdots \end{pmatrix}.$$

$\mathbf{X}_{\text{norte}} \qquad \mathbf{g}_{\text{norte}}(\mathbf{X})$

por un fijo \mathbf{X} , dejar $\mathbf{z} = \mathbf{g}(\mathbf{X})$. Suponer que \mathbf{g} es invertible; por eso $\mathbf{X} = \mathbf{g}^{-1}(\mathbf{z})$. Cualquier infinitesimal n -rectángulo dimensional en \mathbf{X} con volumen V se transforma en un n -paralelepípedo dimensional en \mathbf{z} con volumen $V / |j_{\mathbf{X}}(\mathbf{g})|$, donde $j_{\mathbf{X}}(\mathbf{g})$ es la matriz de jacobiano de la transformación \mathbf{g} , eso es,

$$j_{\mathbf{X}}(\mathbf{g}) = \begin{pmatrix} \frac{\partial g_1}{\partial x_1} & \cdots & \frac{\partial g_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial x_1} & \cdots & \frac{\partial g_n}{\partial x_n} \end{pmatrix}.$$

Ahora considere un vector de columna aleatoria $\mathbf{Z} = \mathbf{g}(\mathbf{X})$. Dejar C ser un pequeño cubo alrededor de \mathbf{z} con volumen h_{norte} . Dejar D ser la imagen de C por debajo de \mathbf{g}^{-1} . Entonces, como en el caso lineal,

$$\text{PAGS}(\mathbf{Z} \in C) \approx h_{\text{norte}} f_{\mathbf{Z}}(\mathbf{z}) \approx h_{\text{norte}} / |j_{\mathbf{X}}(\mathbf{g}^{-1})| f_{\mathbf{X}}(\mathbf{X}).$$

Por lo tanto tenemos la regla de transformación

$$f_{\mathbf{Z}}(\mathbf{z}) = f_{\mathbf{X}}(\mathbf{g}^{-1}(\mathbf{z})) / |j_{\mathbf{X}}(\mathbf{g}^{-1})|, \quad \mathbf{z} \in \mathbb{R}_{\text{norte}}. \quad (1.20)$$

(Nota: $|j_{\mathbf{X}}(\mathbf{g}^{-1})| = 1 / |j_{\mathbf{X}}(\mathbf{g})|$.)

Observación 1.8.1 En la mayoría de las transformaciones de coordenadas, es \mathbf{g}^{-1} que se da, es decir, una expresión para \mathbf{X} como una función de \mathbf{z} más bien que \mathbf{g} .

1.9 TRANSFORMACIONES

Muchos cálculos y manipulaciones que implican distribuciones de probabilidad se facilitan mediante el uso de transformadas. Dos ejemplos típicos son los *función generadora de probabilidad* de una variable aleatoria de valor entero positivo $norte$, definido por

$$GRAMO(z) = MI[z^{norte}] = \sum_{k=0}^{\infty} z^k P(norte=k), \quad |z| < 1,$$

y el *Transformada de Laplace* de una variable aleatoria positiva X definido, por $s > 0$, por

$$L(s) = MI[e^{-sX}] = \begin{cases} \sum_{k=0}^{\infty} e^{-sk} P(X=k) & \text{caso discreto,} \\ \int_0^{\infty} e^{-sx} f(x) dx & \text{caso continuo.} \end{cases}$$

Todas las transformaciones comparten una importante *propiedad de unicidad*: dos distribuciones son iguales si y solo si sus respectivas transformadas son iguales.

■ EJEMPLO 1.8

Dejar $METRO \sim \text{Poi}(\mu)$; entonces su función generadora de probabilidad está dada por

$$GRAMO(z) = \sum_{k=0}^{\infty} z^k \frac{e^{-\mu} \mu^k}{k!} = e^{-\mu} \sum_{k=0}^{\infty} \frac{(z\mu)^k}{k!} = e^{-\mu} e^{z\mu} = e^{-\mu(1-z)}. \quad (1.21)$$

Ahora dejar $norte \sim \text{Poi}(\nu)$ independientemente de $METRO$. Entonces la función generadora de probabilidad de $METRO + norte$ dado por

$$MI[Z_{METRO+norte}] = MI[Z_{METRO}]MI[Z_{norte}] = e^{-\mu(1-z)} e^{-\nu(1-z)} = e^{-(\mu+\nu)(1-z)}.$$

Así, por la propiedad de unicidad, $METRO + norte \sim \text{Poi}(\mu + \nu)$.

■ EJEMPLO 1.9

La transformada de Laplace de $X \sim \text{Gama}(\alpha, \lambda)$ es dado por

$$\begin{aligned} MI[e^{-sX}] &= \int_0^{\infty} e^{-sx} \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)} dx \\ &= \frac{\lambda^\alpha}{(\lambda+s)^\alpha} \int_0^{\infty} e^{-t} t^{\alpha-1} dt \\ &= \frac{\lambda^\alpha}{(\lambda+s)^\alpha}. \end{aligned}$$

Como caso especial, la transformada de Laplace de la $\text{Exp}(\lambda)$ la distribución está dada por $\lambda/(\lambda+s)$. Ahora deja X_1, \dots, X_{norte} ser iid $\text{Exp}(\lambda)$ variables aleatorias. La transformada de Laplace de $S_{norte} = X_1 + \dots + X_{norte}$ es

$$MI[e^{-sS_{norte}}] = MI[e^{-sX_1} \dots e^{-sX_{norte}}] = MI[e^{-sX_1}] \dots MI[e^{-sX_{norte}}] = \left(\frac{\lambda}{\lambda+s} \right)^{norte},$$

lo que demuestra que $S_{norte} \sim \text{Gama}(norte, \lambda)$.

1.10 VARIABLES ALEATORIAS NORMALES CONJUNTAS

Es útil ver las variables aleatorias normalmente distribuidas como simples transformaciones de *normal estándar*—eso es, $N(0,1)$ -distribuido — variables aleatorias. particular, deja $X \sim N(0,1)$. Después X tiene densidad f_X dada por

$$f_X(x) = \sqrt{\frac{1}{2\pi}} e^{-\frac{x^2}{2}}.$$

Ahora considere la transformación $Z = m + \sigma X$. Entonces, por (1.15), Z tiene densidad

$$f_Z(z) = \sqrt{\frac{1}{2\pi\sigma^2}} e^{-\frac{(z-\mu)^2}{2\sigma^2}}.$$

En otras palabras, $Z \sim N(\mu, \sigma^2)$. También podemos expresar esto de la siguiente manera: si $Z \sim N(\mu, \sigma^2)$, después $(Z - \mu)/\sigma \sim N(0,1)$. Este procedimiento se llama *Estandarización*.

Ahora generalizamos esto a *n* dimensiones. Dejar X_1, \dots, X_n ser variables aleatorias normales estándar e independientes. El pdf conjunto de $\mathbf{X} = (X_1, \dots, X_n)$ es dado

$$f_{\mathbf{X}}(\mathbf{x}) = (2\pi)^{-n/2} e^{-\frac{1}{2}\mathbf{x}^T\mathbf{x}}, \mathbf{x} \in \mathbb{R}^n. \quad (1.22)$$

Considera la transformación (es decir, una transformación lineal más un vector constante)

$$\mathbf{Z} = \mathbf{m} + \mathbf{B}\mathbf{X} \quad (1.23)$$

para algunos $m \times n$ matriz \mathbf{B} . Nótese que, por el Teorema 1.8.1, \mathbf{Z} tiene expectativa vector \mathbf{m} y matriz de covarianza $\Sigma = \mathbf{B}\mathbf{B}^T$. Cualquier vector aleatorio de la forma (1.23) se dice que tiene un *conjuntamente normales* o *normal multivariado* distribución. Nosotros escribimos $\mathbf{Z} \sim N(\mathbf{m}, \Sigma)$. Suponer que \mathbf{B} es un invertible $n \times n$ matriz. $\mathbf{Y} = \mathbf{Z} - \mathbf{m}$ entonces, por (1.19), la densidad de

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{|\mathbf{B}| (2\pi)^{n/2}} e^{-\frac{1}{2}\mathbf{y}^T \mathbf{B}^{-1} \mathbf{B}^{-T} \mathbf{y}} = \frac{1}{|\mathbf{B}| (2\pi)^{n/2}} e^{-\frac{1}{2}\mathbf{y}^T \Sigma^{-1} \mathbf{y}}.$$

Tenemos $|\mathbf{B}| = |\Sigma|^{1/2}$ (como $\mathbf{y}^T \mathbf{B}^{-1} \mathbf{B}^{-T} \mathbf{y} = \mathbf{y}^T \Sigma^{-1} \mathbf{y}$), de modo que

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}\mathbf{y}^T \Sigma^{-1} \mathbf{y}}.$$

Porque \mathbf{Z} se obtiene de \mathbf{Y} simplemente agregando un vector constante \mathbf{m} , tenemos $f_{\mathbf{Z}}(\mathbf{z}) = f_{\mathbf{Y}}(\mathbf{z} - \mathbf{m})$, y por lo tanto

$$f_{\mathbf{Z}}(\mathbf{z}) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{z} - \mathbf{m})^T \Sigma^{-1} (\mathbf{z} - \mathbf{m})}, \mathbf{z} \in \mathbb{R}^n. \quad (1.24)$$

Tenga en cuenta que esta fórmula es muy similar a la del caso unidimensional.

Por el contrario, dada una matriz de covarianza $\Sigma = (\sigma_{ij})$, existe una única matriz triangular inferior

$$\mathbf{B} = \begin{pmatrix} b_{11} & 0 & \cdots & 0 \\ b_{21} & b_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \cdots & b_{nn} \end{pmatrix} \quad (1.25)$$

tal que $\Sigma = \mathbf{B}\mathbf{B}^T$. Esta matriz se puede obtener de manera eficiente a través de la *Método de la raíz cuadrada de Cholesky*; véase la Sección A.1 del Apéndice.

1.11 TEOREMAS DEL LÍMITE

Discutimos brevemente dos de los principales resultados en probabilidad: la ley de los grandes números y el teorema del límite central. Ambos están asociados con sumas de variables aleatorias independientes.

Dejar X_1, X_2, \dots ser iid variables aleatorias con expectativa μ y varianza σ^2 . Para cada n , dejar $S_n = X_1 + \dots + X_n$. Ya que X_1, X_2, \dots son iid, tenemos $MI[S_n] = nMI[X_1] = n\mu$ y $Var(S_n) = nVar(X_1) = n\sigma^2$.

La ley de los grandes números establece que S_n/n está cerca de μ para n grande. Aquí está la declaración más precisa.

Teorema 1.11.1 (Ley fuerte de los grandes números) Si X_1, \dots, X_n son iid con expectativa μ , entonces

$$\lim_{n \rightarrow \infty} \frac{S_n}{n} = \mu \text{ casi seguramente.}$$

El teorema del límite central describe la distribución límite de S_n/n (o S_n/n), y se aplica tanto a variables aleatorias continuas como discretas. En términos generales, establece que la suma aleatoria S_n tiene una distribución que es aproximadamente normal, cuando n es largo. La declaración más precisa se da a continuación.

Teorema 1.11.2 (Teorema del límite central) Si X_1, \dots, X_n son iid con expectativa μ y varianza $\sigma^2 < \infty$, entonces para todos $x \in \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \frac{S_n - n\mu}{\sigma\sqrt{n}} = \Phi(x),$$

dónde Φ es la cdf de la distribución normal estándar.

En otras palabras, S_n tiene una distribución que es aproximadamente normal, con expectativa $n\mu$ y varianza $n\sigma^2$. Para ver el teorema del límite central en acción, considere la figura 1.2. La parte izquierda muestra los pdfs de S_1, \dots, S_4 para el caso en que el $\{X_i\}$ tener una $[0,1]$ distribución. La parte derecha muestra lo mismo para el $\text{Exp}(1)$ distribución. Vemos claramente la convergencia a una curva en forma de campana, característica de la distribución normal.

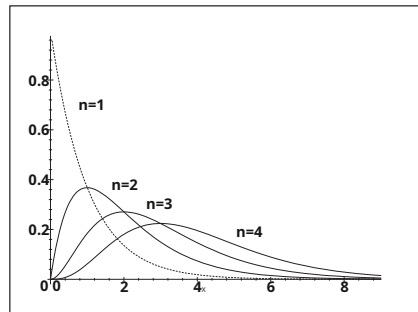
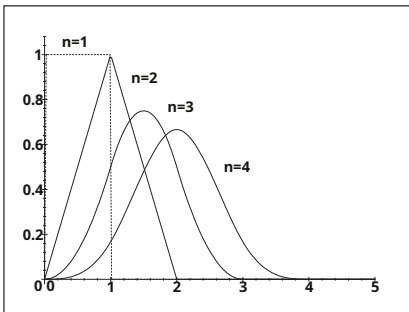


Figura 1.2: Ilustración del teorema del límite central para (izquierda) la distribución uniforme y (derecha) la distribución exponencial.

Una consecuencia directa del teorema del límite central y del hecho de que un $\text{Compartimiento}(norte, pag)$ variable aleatoria X puede verse como la suma de $norte$ iid $\text{Ber}(pag)$ variables aleatorias, $X = X_1 + \dots + X_{norte}$, es que para grandes $norte$

$$\text{PAGS}(X-k) \approx \text{PAGS}(Y-k), \quad (1.26)$$

con $Y \sim \text{NORTE}(np, np(1-pag))$. Como regla general, esta *aproximación normal a la distribución binomial* es exacto si ambos $norte$ y $norte(1-pag)$ son mayores que 5.

También hay un teorema del límite central para vectores aleatorios. La versión multidimensional es la siguiente: Sea X_1, \dots, X_{norte} ser iid vectores aleatorios con vector de expectativa μ y matriz de covarianza Σ . Entonces para grandes $norte$ el vector aleatorio $X_1 + \dots + X_{norte}$ tiene aproximadamente una distribución normal multivariante con vector de expectativa $norte\mu$ y matriz de covarianza $norte\Sigma$.

1.12 PROCESOS VENENOSOS

El proceso de Poisson se usa para modelar ciertos tipos de llegadas o patrones. Imagine, por ejemplo, un telescopio que pueda detectar fotones individuales de una galaxia lejana. Los fotones llegan en momentos aleatorios. T_1, T_2, \dots Dejar $norte_t$ denotar el número de llegadas en el intervalo de tiempo $[0, t]$, eso es, $norte_t = \text{sorber}\{k: T_k \leq t\}$. Tenga en cuenta que el número de llegadas en un intervalo $y = (un, b]$ es dado por $norte_b - norte_a$. También lo denotaremos por $norte(un, b]$. Una ruta de muestra del proceso de conteo de llegadas $\{norte_t, t \geq 0\}$ se da en la figura 1.3.

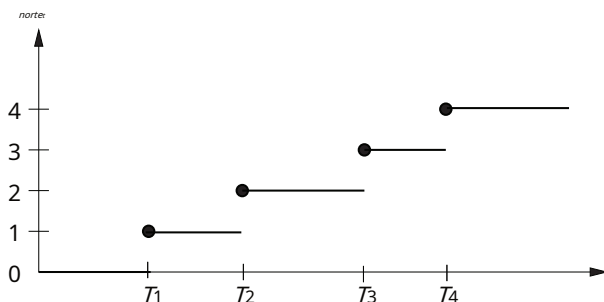


Figura 1.3: Una ruta de muestra del proceso de conteo de llegadas $\{norte_t, t \geq 0\}$.

Para este proceso de llegada en particular, se supondría que el número de llegadas en un intervalo $(un, b]$ es independiente del número de llegadas en el intervalo $(discos compactos)$ cuando los dos intervalos no se cortan. Tales consideraciones conducen a la siguiente definición:

Definición 1.12.1 (Proceso de Poisson) Un proceso de conteo de llegadas $norte_t = \{NORTE_t; t \geq 0\}$ se llama un *proceso de envenenamiento* con tasa $\lambda > 0$ si

- Los números de puntos en intervalos que no se superponen son independientes.
- El número de puntos en el intervalo y tiene una distribución de Poisson con media $\lambda \times \text{longitud}(y)$.

Combinando (a) y (b) vemos que el número de llegadas en cualquier intervalo pequeño $(t, t+h]$ es independiente del proceso de llegada hasta el momento t y tiene un $\text{Poi}(\lambda h)$ distribución. En particular, la probabilidad condicional de que ocurra exactamente una llegada durante el intervalo de tiempo $(t, t+h]$ es $\text{PAGS}(\text{norte}(t, t+h] = 1/N_t) = \text{mi}_{-1} \lambda h \lambda h \approx \lambda h$. De manera similar, la probabilidad de que no haya llegadas es de aproximadamente $1 - \lambda h$. Para pequeños h . En otras palabras, λ es la *Velocidad* en la que se producen las llegadas. Nótese también que desde $\text{norte} \sim \text{Poi}(\lambda t)$, el número esperado de llegadas en $[0, t]$ es λt , eso es, $\text{MI}[\text{norte}] = \lambda t$. En la Definición 1.12.1 norte se ve como una medida de conteo aleatorio, donde $\text{norte}(y)$ cuenta el número aleatorio de llegadas en el conjunto y .

Una relación importante entre norte y T_{norte} es

$$\{\text{NORTE} \leq \text{norte}\} = \{T_{\text{norte}} \leq t\}. \quad (1.27)$$

En otras palabras, el número de llegadas en $[0, t]$ Por lo menos norte si y solo si el norte -la llegada se produce a la hora o antes t . Como consecuencia, tenemos

$$\begin{aligned} \text{PAGS}(T_{\text{norte}} \leq t) &= \text{PAGS}(\text{norte} \leq \text{norte}) = 1 - \sum_{k=0}^{\infty} \text{PAGS}(\text{norte} = k) \\ &= 1 - \sum_{k=0}^{\infty} \text{mi}_{-1} \lambda^k \frac{(\lambda t)^k}{k!}, \end{aligned}$$

que corresponde exactamente a la cdf de la $\text{Gama}(\text{norte}, \lambda)$ distribución; vea el problema 1.17. De este modo

$$T_{\text{norte}} \sim \text{Gama}(\text{norte}, \lambda). \quad (1.28)$$

Por lo tanto, cada T_{norte} tiene la misma distribución que la suma de norte independientes $\text{Exp}(\lambda)$ -variables aleatorias distribuidas. Esto se corresponde con la segunda caracterización importante de un proceso de Poisson:

Un proceso de conteo de llegadas $\{\text{norte}_t\}$ es un proceso de Poisson con tasa λ si y solo si los tiempos entre llegadas $A_1 = T_1, A_2 = T_2 - T_1, \dots$ son independientes y $\text{Exp}(\lambda)$ -Variables aleatorias distribuidas.

Los procesos de Poisson y Bernoulli son similares, y se puede aprender mucho sobre los procesos de Poisson a través de lo siguiente *Aproximación de Bernoulli*. Dejar $\text{norte} = \{\text{NORTE}_t\}$ sea un proceso de Poisson con parámetro λ . Dividimos el eje de tiempo en pequeños intervalos de tiempo $[0, h), [h, 2h), \dots$ y cuente cuántas llegadas ocurren en cada intervalo. Tenga en cuenta que el número de llegadas en cualquier pequeño intervalo de tiempo de longitud h es, con alta probabilidad, 1 (con probabilidad $\lambda h \text{mi}_{-1} \lambda h \approx \lambda h$) o 0 (con probabilidad $e^{-\lambda h} \approx 1 - \lambda h$). A continuación, defina $X = \{X_{\text{norte}}\}$ ser un proceso de Bernoulli con parámetro de éxito $\text{pags} = \lambda h$. Poner $Y_0 = 0$ y deja $Y_{\text{norte}} = X_1 + \dots + X_{\text{norte}}$ Sea el número total de éxitos en norte juicios $Y = \{Y_{\text{norte}}\}$ se llama el *Aproximación de Bernoulli a norte*. podemos ver norte como un caso límite de Y a medida que disminuimos h .

Como ejemplo de la utilidad de esta interpretación, ahora demostramos que la propiedad de Poisson (b) en la Definición 1.12.1 se sigue básicamente de la *independencia* supuesto (a). Para pequeños h , norte debe tener aproximadamente la misma distribución

como Y_{norte} , dónde $norte$ es la parte entera de t/h (nosotros escribimos $norte = t/h$). Por eso,

$$\begin{aligned} \text{PAGS}(norte = k) &\approx \text{PAGS}(Y_{norte} = k) \\ &= \binom{norte}{k} (\lambda h)^k (1 - (\lambda h))^{norte-k} \\ &\approx \binom{norte}{k} (\lambda t/n)^k (1 - (\lambda t/n))^{norte-k} \\ &\approx \frac{mi^{-\lambda} t (\lambda t)^k}{k!}. \end{aligned} \quad (1.29)$$

La ecuación (1.29) se deriva de la aproximación de Poisson a la distribución binomial; vea el problema 1.22.

Otra aplicación de la aproximación de Bernoulli es la siguiente. Para el proceso de Bernoulli, dado que el número total de éxitos es k , las posiciones de los k éxitos se distribuyen uniformemente en los puntos $1, \dots, norte$. La propiedad correspondiente para el proceso de Poisson $norte$ es eso dado $norte = norte$, los tiempos de llegada T_1, \dots, T_{norte} se distribuyen según las estadísticas del pedido $X_{(1)}, \dots, X_{(norte)}$, dónde X_1, \dots, X_{norte} son iid $u[0, t]$.

1.13 PROCESOS DE MARKOV

Los procesos de Markov son procesos estocásticos cuyos futuros son condicionalmente independientes de sus pasados dados sus valores presentes. Más formalmente, un proceso estocástico $\{X_t, t \in T\}$, con $T \subseteq \mathbb{R}$, se llama un *proceso de Markov* si por cada $s > 0$ y t ,

$$(X_{t+s} / X_t, t \leq t) \sim (X_{t+s} / X_t). \quad (1.30)$$

En otras palabras, la distribución condicional de la variable futura X_{t+s} , dado todo el pasado del proceso $\{X_t, t \leq t\}$, es lo mismo que la distribución condicional de X_{t+s} dado solo el presente X_t . Es decir, para predecir estados futuros, solo necesitamos conocer el presente. La propiedad (1.30) se llama *Propiedad de Markov*.

Dependiendo del conjunto de índices T y espacio de estado m (el conjunto de todos los valores $\{X_t\}$ puede tomar), los procesos de Markov vienen en muchas formas diferentes. Un proceso de Markov con un conjunto de índices discretos se llama *cadena de Markov*. Un proceso de Markov con un espacio de estado discreto y un conjunto de índices continuo (como \mathbb{R}_+) se llama *Proceso de salto de Markov*.

1.13.1 Cadenas de Markov

Considere una cadena de Markov $X = \{X_t, t \in norte\}$ con un espacio de estado discreto (es decir, contable) m . En este caso la propiedad de Markov (1.30) es

$$\text{PAGS}(X_{t+1} = X_{t+1} / X_0 = X_0, \dots, X_t = X_t) = \text{PAGS}(X_{t+1} = X_{t+1} / X_t = X_t) \quad (1.31)$$

para todos $X_0, \dots, X_{t+1}, \in m, y t \in norte$. Nos restringimos a las cadenas de Markov para las cuales las probabilidades condicionales

$$\text{PAGS}(X_{t+1} = j / X_t = i), \text{ y } o, j \in m \quad (1.32)$$

son independientes del tiempo t . Tales cadenas se llaman *homogéneas en el tiempo*. Las probabilidades en (1.32) se llaman *(un paso) probabilidades de transición* de X . La distribución de X_0 se llama *distribución inicial* de la cadena de Markov. Las probabilidades de transición de un paso y la distribución inicial especifican completamente la distribución de X . A saber, tenemos por la regla del producto (1.4) y la propiedad de Markov (1.30),

$$\begin{aligned} \text{PAGS}(X_0=X_0, \dots, X_t=X_t) \\ = \text{PAGS}(X_0=X_0) \text{PAGS}(X_1=X_1/X_0=X_0) \cdots \text{PAGS}(X_t=X_t/X_0=X_0, \dots, X_{t-1}=X_{t-1}) = \text{PAGS}(X_0=X_0) \\ \text{PAGS}(X_1=X_1/X_0=X_0) \cdots \text{PAGS}(X_t=X_t/X_{t-1}=X_{t-1}). \end{aligned}$$

Ya que m es contable, podemos organizar las probabilidades de transición de un paso en una matriz. Esta matriz se llama *(un paso) matriz de transición* de X . Normalmente denotamos por PAGS . por ejemplo, cuando $m=\{0,1,2,\dots\}$, la matriz de transición PAGS tiene la forma

$$\text{PAGS} = \begin{pmatrix} \text{pags}_{00} & \text{pags}_{01} & \text{pags}_{02} & \dots \\ \text{pags}_{10} & \text{pags}_{11} & \text{pags}_{12} & \dots \\ \text{pags}_{20} & \text{pags}_{21} & \text{pags}_{22} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

Tenga en cuenta que los elementos en cada fila son positivos y suman la unidad.

Otra forma conveniente de describir una cadena de Markov X es a través de su *gráfico de transición*. Los estados se indican mediante los nodos del grafo y un signo estrictamente positivo (>0) probabilidad de transición $\text{pags}_{y \rightarrow o}$ del estado i a j se indica con una flecha desde i a j con peso $\text{pags}_{y \rightarrow o}$.

■ EJEMPLO 1.10 Paseo aleatorio sobre los enteros

Dejar pags ser un número entre 0 y 1. La cadena de Markov X con espacio de estado \mathbb{Z} y matriz de transición PAGS definido por

$$\text{PAGS}(y \rightarrow y+1) = \text{pags} \quad \text{pags}(y \rightarrow y-1) = q = 1 - \text{pags}, \quad \text{para todos } i \in \mathbb{Z}$$

se llama un *paseo aleatorio en los enteros*. Dejar X empezar en 0; de este modo, $\text{PAGS}(X_0=0) = 1$. El gráfico de transición correspondiente se muestra en la Figura 1.4. Comenzando en 0, la cadena toma pasos subsiguientes hacia la derecha con probabilidad pags y a la izquierda con probabilidad q .

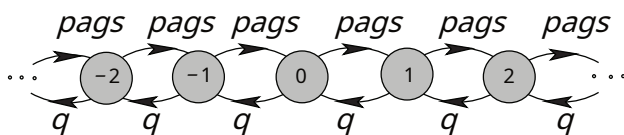


Figura 1.4: Gráfico de transición para un paseo aleatorio en \mathbb{Z}

Mostramos a continuación cómo calcular la probabilidad de que, partiendo del estado i en algún tiempo (discreto) t , estamos en j en tiempo (discreto) $t+s$, es decir, la probabilidad $\text{PAGS}(X_{t+s}=j/X_t=i)$. Para mayor claridad, supongamos que $m=\{1,2,\dots,m\}$ para algunos fijos m , de modo que PAGS es un $m \times m$ matriz. Para $t=0,1,2,\dots$, define el vector fila

$$\pi_t = (\text{PAGS}(X_t=1), \dots, \text{PAGS}(X_t=m)).$$

Nosotros llamamos $\pi_{(0)}$ a *vectores de distribución*, o simplemente *la distribución*, de X en el momento t y $\pi_{(0)}$ a *la distribución inicial* de X . El siguiente resultado muestra que las probabilidades de paso se pueden encontrar simplemente mediante la multiplicación de matrices.

Teorema 1.13.1 *La distribución de X en el momento t es dado por*

$$\pi_{(t)} = \pi_{(0)} \text{PAGS}_t \quad (1.33)$$

para todo $t=0, 1, \dots$ (Aquí PAGS_0 denota la matriz identidad.)

Prueba: La prueba es por inducción. La igualdad (1.33) se cumple para $t=0$ por definición. Supongamos que esta igualdad es cierta para algunos $t=0, 1, \dots$. Tenemos

$$\text{PAGS}(X_{t+1}=k) = \sum_{j=1}^n \text{PAGS}(X_{t+1}=k / X_t=j) \text{PAGS}(X_t=j).$$

Pero se supone que (1.33) es cierta para t , así que $\text{PAGS}(X_t=j)$ es el j -ésimo elemento de $\pi_{(0)} \text{PAGS}_t$. Es más, $\text{PAGS}(X_{t+1}=k / X_t=j)$ es el (y_0, k) -ésimo elemento de PAGS . Por lo tanto, por cada k ,

$$\sum_{j=1}^n \text{PAGS}(X_{t+1}=k / X_t=j) \text{PAGS}(X_t=j) = \sum_{j=1}^n \text{PAGS}(y_0, k) (\pi_{(0)} \text{PAGS}_t)(j),$$

que es solo el k -ésimo elemento de $\pi_{(0)} \text{PAGS}_{t+1}$. Esto completa el paso de inducción y, por lo tanto, se prueba el teorema. \square

Tomando $\pi_{(0)}$ como el i -ésimo vector unitario, \mathbf{m}_i , las probabilidades de transición de paso se pueden encontrar como $\text{PAGS}(X_t=j / X_0=i) = (\mathbf{m}_i \text{PAGS}_t)(j) = \text{PAGS}_t(y_0, j)$. Cuál es el (y_0, j) -ésimo elemento de matriz PAGS_t . Así, para encontrar las probabilidades de transición de paso, solo tenemos que calcular el t -ésima potencia de PAGS .

1.13.2 Clasificación de Estados

Dejar X ser una cadena de Markov con espacio de estado discreto \mathbf{m} y matriz de transición PAGS . Podemos caracterizar las relaciones entre estados de la siguiente manera: Si los estados i, j son tales que $\text{PAGS}_t(y_0, j) > 0$ para algunos $t \geq 0$, decimos que *i lleva a j* escribiendo $i \rightarrow j$. Nosotros decimos eso *i, j comunicarse* si $i \rightarrow j$ y $j \rightarrow i$, y escribiendo $i \leftrightarrow j$. Usando la relación " \leftrightarrow ", podemos dividir \mathbf{m} dentro de *clases de equivalencia* tales que todos los estados en una clase de equivalencia se comunican entre sí pero no con ningún estado fuera de esa clase. Si solo hay una clase de equivalencia, decimos que la cadena de Markov es *irreducible*. Si un conjunto de estados A es tal que $\sum_{j \in A} \text{PAGS}_t(y_0, j) = 1$ para todos $i \in A$, después A se llama un *cerrado* establecer. Un estado i se llama un *absorbente* si $\{i\}$ está cerrado. Por ejemplo, en el gráfico de transición representado en la Figura 1.5, las clases de equivalencia son $\{1, 2\}$, $\{3\}$, y $\{4, 5\}$. Clase $\{1, 2\}$ es el único conjunto cerrado: la cadena de Markov no puede escapar de él. Si faltara el estado 1, el estado 2 sería absorbente. En el ejemplo 1.10, la cadena de Markov es irreducible ya que todos los estados se comunican.

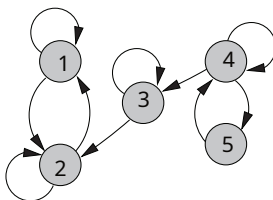


Figura 1.5: Un gráfico de transición con tres clases de equivalencia.

Otra clasificación de estados se obtiene observando el sistema desde un punto de vista local. En particular, deja T_i indicar el momento en que la cadena visita por primera vez el estado i , o primero vuelve a i si empezó ahí, y deja n denote el número total de visitas a j desde el tiempo 0 en adelante. Nosotros escribimos $PAGS(A)$ por $PAGS(A/X_0=j)$ para cualquier evento A . Denotamos el operador de expectativa correspondiente por m_i . Estado j se llama un *recurrente* indicar si $PAGS(T_j < \infty) = 1$; de lo contrario, j se llama *transitorio*. Un estado recurrente se llama *positivo recurrente* si $m_i[T_j < \infty] < \infty$; de lo contrario, se llama *nulo recurrente*. Finalmente, se dice que un estado *esperiódico*, con período δ , si $\delta \geq 2$ es el entero más grande para el cual $PAGS(T_j = n\delta \text{ para algunos } n) = 1$; de lo contrario, se llama *aperiódico*. Por ejemplo, en la figura 1.5, los estados 1 y 2 son recurrentes y los otros estados son transitorios. Todos estos estados son aperiódicos. Los estados de la caminata aleatoria del ejemplo 1.10 son periódicos con período 2.

Se puede demostrar que la recurrencia y la transitoriedad son propiedades de clase. En particular, si $i \leftrightarrow j$, después i recurrente (transitorio) $\Leftrightarrow j$ recurrente (transitorio). Por lo tanto, en una cadena de Markov irreducible, que un estado sea recurrente implica que todos los demás estados también son recurrentes. Y si un estado es transitorio, también lo son todos los demás.

1.13.3 Comportamiento limitante

El comportamiento límite o de "estado estacionario" de las cadenas de Markov como $t \rightarrow \infty$ es de considerable interés e importancia, y este tipo de comportamiento es a menudo más simple de describir y analizar que el comportamiento "transitorio" de la cadena para fijos t . Se puede demostrar (ver, por ejemplo, [3]) que en una cadena de Markov irreducible y aperiódica con matriz de transición $PAGS$ la Las probabilidades de paso convergen a una constante que no depende del estado inicial. Más específicamente,

$$\lim_{t \rightarrow \infty} PAGS(y_0, j) = \pi_j \quad (1.34)$$

por algún número $0 < \pi_j < 1$. Además, $\pi_j > 0$ si j es recurrente positivo y $\pi_j = 0$ de lo contrario. La razón intuitiva detrás de este resultado es que el proceso "olvida" dónde estaba inicialmente si se prolonga lo suficiente. Esto es cierto tanto para las cadenas de Markov finitas como para las numerables infinitas. Los números $\{\pi_j, j \in M\}$ formar el *limitando distribución* de la cadena de Markov, siempre que $\pi_j > 0$ y $\sum_j \pi_j = 1$. Nota que estas condiciones no siempre se cumplen: claramente no se cumplen si la cadena de Markov es transitoria, y pueden no cumplirse si la cadena de Markov es recurrente (es decir, cuando los estados son nulos-recurrentes). El siguiente teorema da un método para obtener distribuciones límite. Aquí suponemos por simplicidad que $M = \{0, 1, 2, \dots\}$. La distribución límite se identifica con el vector fila $\pi = (\pi_0, \pi_1, \dots)$.

Teorema 1.13.2 Para una cadena de Markov irreducible y aperiódica con matriz de transición $PAGS$, si la distribución límite π existe, entonces está determinada únicamente por la solución de

$$\pi = \pi PAGS, \quad (1.35)$$

con $\pi_j \geq 0$ y $\sum_j \pi_j = 1$. Por el contrario, si existe un vector fila positivo π satisfactorio a (1.35) y sumando 1, entonces π es la distribución límite de la cadena de Markov. Además, en ese caso, $\pi_j > 0$ para todos j y todos los estados son recurrentes positivos.

Prueba. (Sketch) Para el caso en que m es finito, el resultado es simplemente una consecuencia de (1.33). Es decir, con $\pi_{(0)}$ siendo el i -ésimo vector unitario, tenemos

$$PAGS_{t+1}(y_0, j) = \pi_{(0)} PAGS_t PAGS(j) = \sum_{k \in m} PAGS(y_0, k) PAGS(k, j).$$

Alquiere $t \rightarrow \infty$, obtenemos (1.35) de (1.34), siempre que podamos cambiar el orden del límite y el \sum suma. Para mostrar la unicidad, suponga que otro vector y , con $y_j \geq 0$ y $\sum_j y_j = 1$, satisface $y = y PAGS$. Entonces es fácil mostrar por inducción que $y = y PAGS_t$ para cada t . Por lo tanto, dejar $t \rightarrow \infty$, obtenemos para cada j

$$y_j = \sum_i y_i \pi_{ij} = \pi_j,$$

desde el $\{y_j\}$ suma a la unidad. Omitimos la demostración del enunciado inverso. -

■ EJEMPLO 1.11 Paseo aleatorio sobre los enteros positivos

Esta es una caminata aleatoria ligeramente diferente a la del Ejemplo 1.10. Dejar X ser un paseo al azar en $m = \{0, 1, 2, \dots\}$ con matriz de transición

$$PAGS = \begin{pmatrix} q & pag & 0 & \dots \\ 0 & q & pag & 0 & \dots \\ 0 & q & 0 & pag & \dots \\ \dots & \dots & \dots & \dots & \dots \end{pmatrix},$$

donde $0 < pag < 1$ y $q = 1 - pag$. X_t podría representar, por ejemplo, el número de clientes que están esperando en una cola a la vez t .

Todos los estados pueden alcanzarse entre sí, por lo que la cadena es irreducible y cada estado es recurrente o transitorio. La ecuación $\pi = \pi PAGS$ se convierte

$$\begin{aligned} \pi_0 &= q \pi_0 + q \pi_1, \\ \pi_1 &= pag \pi_0 + q \pi_2, \\ \pi_2 &= pag \pi_1 + q \pi_3, \\ \pi_3 &= pag \pi_2 + q \pi_4, \end{aligned}$$

y así. Podemos resolver este conjunto de ecuaciones secuencialmente. si dejamos $r = p/q$, entonces podemos expresar la π_1, π_2, \dots en términos de π_0 y r como

$$\pi_j = r^j \pi_0, \quad j = 0, 1, 2, \dots$$

Si $p < q$, después $r < 1$ un $\sum_{j=0}^{\infty} \pi_j = \pi_0(1-r)$, y al elegir $\pi_0 = 1-r$, $\pi_j = 1$. Por lo tanto, podemos hacer la suma para $r < 1$, hemos encontrado la limitante distribución $\pi = (1-r)(1, r, r^2, r^3, \dots)$ para este Markov chain. Los estados son recurrentes, y todo el por lo tanto positivos recurrentes. Sin embargo cuando $p = q$, $\pi_j = 0$ o infinito, y por lo tanto todos los estados son nulos-recurrentes o transitorios. (Se puede demostrar que sólo el caso $p = q$ conduce a estados recurrentes nulos.)

Dejar X sea una cadena de Markov con distribución limitante π . Suponer $\pi_0 = \pi$. Entonces, combinando (1.33) y (1.35), tenemos $\pi_t = \pi$. Por tanto, si la distribución inicial de la cadena de Markov es igual a la distribución límite, entonces la distribución de X_t es igual para todos t (y está dada por esta distribución límite). De hecho, no es difícil demostrar que para cualquier k la distribución de $X_k, X_{k+1}, X_{k+2}, \dots$ es el mismo que el de X_0, X_1, \dots . En otras palabras, cuando $\pi_0 = \pi$, la cadena de Markov es un proceso estocástico estacionario. Más formalmente, un proceso estocástico $\{X_t, t \in \mathbb{N}\}$ se llama *estacionario* si, para cualquier positivo t, t_1, \dots, t_n , el vector $(X_t, \dots, X_{t+t_1}, \dots, X_{t+t_1+t_2}, \dots, X_{t+t_1+t_2+t_3}, \dots, X_{t+t_1+t_2+t_3+t_4}, \dots)$ tiene la misma distribución que $(X_0, \dots, X_{t_1}, \dots, X_{t_1+t_2}, \dots, X_{t_1+t_2+t_3}, \dots, X_{t_1+t_2+t_3+t_4}, \dots)$. Definiciones similares son válidas cuando el conjunto de índices es \mathbb{Z} , \mathbb{R} , o \mathbb{R}^+ . Por esta razón cualquier distribución π para la cual (1.35) se cumple se llama

estación π

distribución π .

Señalando que $\sum_j p_{js} = 1$, podemos reescribir (1.35) como el sistema de ecuaciones

$$\sum_j \pi_j p_{js} = \sum_j \pi_j p_{js} \quad \text{para todos } s \in M \quad (1.36)$$

Estos se llaman *ecuaciones de equilibrio global*. Podemos interpretar (1.35) como el enunciado de que el "flujo de probabilidad" de s se equilibra con el flujo de probabilidad en s . Una generalización importante, que se deriva directamente de (1.36), establece que el mismo equilibrio de los flujos de probabilidad se cumple para un conjunto arbitrario A . Es decir, para cada conjunto A de estados que tenemos

$$\sum_{i \in A} \sum_{j \notin A} \pi_i p_{ji} = \sum_{i \in A} \sum_{j \notin A} \pi_j p_{ji} \quad (1.37)$$

1.13.4 Reversibilidad

La reversibilidad es una noción importante en la teoría de Markov y procesos más generales. Un proceso estocástico estacionario $\{X_t\}$ con juego de índices \mathbb{Z} o \mathbb{R} se ha dicho *reversible* si, para cualquier entero positivo n y para todos t_1, \dots, t_n , el vector $(X_{t_1}, \dots, X_{t_1+t_1}, \dots, X_{t_1+t_1+t_2}, \dots, X_{t_1+t_1+t_2+t_3}, \dots, X_{t_1+t_1+t_2+t_3+t_4}, \dots)$ tiene la misma distribución que $(X_{-t_1}, \dots, X_{-t_1-t_1}, \dots, X_{-t_1-t_1-t_2}, \dots, X_{-t_1-t_1-t_2-t_3}, \dots, X_{-t_1-t_1-t_2-t_3-t_4}, \dots)$. Una forma de visualizar esto es imaginar que hemos tomado un video del proceso estocástico, que podemos ejecutar hacia adelante y hacia atrás. Si no podemos determinar si el video avanza o retrocede, el proceso es reversible. El resultado principal para las cadenas de Markov reversibles es que un proceso de Markov estacionario es reversible si y solo si existe una colección de números positivos $\{\pi_i, i \in M\}$ sumando a la unidad que satisfacen las *ecuaciones de equilibrio detalladas (o locales)*

$$\pi_i p_{ij} = \pi_j p_{ji}, \quad j \in M \quad (1.38)$$

Siempre que tal colección $\{\pi_i\}$ existe, es la distribución estacionaria del proceso.

Una buena manera de pensar en las ecuaciones de equilibrio detalladas es que equilibran el flujo de probabilidad del estado j a estado i con eso del estado i a estado j . Contraste

esto con las ecuaciones de equilibrio (1.36), que equilibran el flujo de probabilidad fuera del estado i con eso en estado i .

criterio de Kolmogorov es un criterio simple para la reversibilidad basado en las probabilidades de transición. Establece que un proceso de Markov estacionario es reversible si y solo si sus tasas de transición satisfacen

$$pags(\hat{h}, \hat{i})pags(\hat{i}, \hat{i}_2) \dots pags(\hat{i}_{n-1}, \hat{i}_{norte})pags(\hat{i}_{norte}, \hat{h}) = pags(\hat{h}, \hat{i}_{norte})pags(\hat{i}_{norte}, \hat{i}_{n-1}) \dots pags(\hat{i}_2, \hat{h}) \quad (1.39)$$

para todos los bucles finitos de estados $\hat{h}, \dots, \hat{i}_{norte}, \hat{h}$. (Para mayor claridad, hemos utilizado la notación $pags(yo, j)$ más bien que $pags_{yo}$ para las probabilidades de transición.) La idea es bastante intuitiva: si es más probable que el proceso en el tiempo hacia adelante atravesase un cierto ciclo cerrado en una dirección que en la dirección opuesta, entonces en el tiempo hacia atrás exhibirá el comportamiento opuesto y, por lo tanto, tener un criterio para detectar la dirección del tiempo. Si tal comportamiento de "bucle" no ocurre, el proceso debe ser reversible.

1.13.5 Procesos de salto de Markov

A *Proceso de salto de Markov* $X = \{X_t, t \geq 0\}$ puede verse como una generalización en tiempo continuo de una cadena de Markov y también de un proceso de Poisson. La propiedad de Markov (1.30) ahora dice

$$PAGS(X_{t+s} = X_{tu} / X_{tu} = X_{tu}, tu - t) = PAGS(X_{t+s} = X_{t+s} / X_t = X_t). \quad (1.40)$$

Como en el caso de la cadena de Markov, por lo general se supone que el proceso es *homogéneo en el tiempo*, eso es, $PAGS(X_{t+s} = j / X_t = i)$ no depende de t . Denote esta probabilidad por $PAGS(yo, j)$. Una cantidad importante es la *tasa de transición* q_{yo} del estado i a j , definido por $i \neq j$ como

$$q_{yo} = \lim_{t \rightarrow 0} \frac{PAGS(yo, j)}{t}.$$

La suma de las tasas fuera del estado i se denota por q_i . Una ruta de muestra típica de X se muestra en la Figura 1.6. El proceso salta a veces T_1, T_2, \dots a los estados Y_1, Y_2, \dots , permaneciendo cierto tiempo en cada estado.

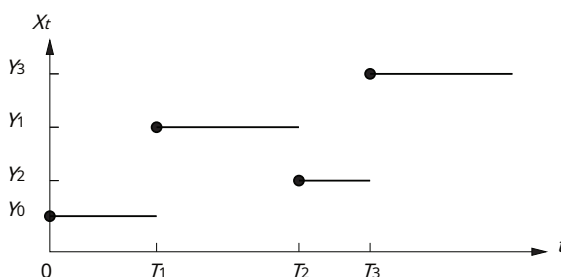


Figura 1.6: Una ruta de muestra de un proceso de salto de Markov $\{X_t, t \geq 0\}$.

Más precisamente, un proceso de salto de Markov X se comporta (bajo condiciones de regularidad adecuadas; ver [3]) de la siguiente manera:

1. Dado su pasado, la probabilidad de que X salte de su estado actual i a estado j es $k_{yo} = q_{yo}/q_i$.

2. La cantidad de tiempo que X gasta en el estado j tiene una distribución exponencial con media $1/q_j$, independiente de su historia pasada.

La primera afirmación implica que el proceso $\{Y_{norte}\}$ es de hecho una cadena de Markov, con matriz de transición $k = (k_{y_0})$.

Una forma conveniente de describir un proceso de salto de Markov es a través de su *gráfico de tasa de transición*. Esto es similar a un gráfico de transición para las cadenas de Markov. Los estados están representados por los nodos del gráfico y una tasa de transición del estado i a j se indica con una flecha desde i a j con peso q_{y_0} .

■ EJEMPLO 1.12 Proceso de nacimiento y muerte

Un *proceso de nacimiento y muerte* es un proceso de salto de Markov con un gráfico de tasa de transición de la forma dada en la Figura 1.7. Imagina que X_t representa el número total de individuos en una población en el tiempo t . Los saltos a la derecha corresponden a nacimientos, y los saltos a la izquierda a defunciones. Las *tasas de natalidad* $\{b_i\}$ y las *tasas de mortalidad* $\{d_i\}$ pueden diferir de un estado a otro. Muchas aplicaciones de las cadenas de Markov involucran procesos de este tipo. Tenga en cuenta que el proceso salta de un estado a

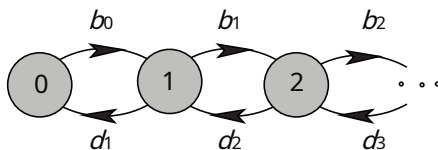


Figura 1.7: El gráfico de tasa de transición de un proceso de nacimiento y muerte.

el siguiente según una cadena de Markov con probabilidades de transición $k_{0,1} = 1$, $k_{y_0, y_0+1} = b_i / (b_i + d_i)$, y $k_{y_0, y_0-1} = d_i / (b_i + d_i)$, $i = 1, 2, \dots$. Además, gasta un $\text{Exp}(b_0)$ cantidad de tiempo en el estado 0 y $\text{Exp}(b_i + d_i)$ en los demás estados.

Comportamiento limitante Ahora formulamos los análogos en tiempo continuo de (1.34) y el Teorema 1.13.2. La irreducibilidad y la recurrencia para los procesos de salto de Markov se definen de la misma manera que para las cadenas de Markov. Por simplicidad, supongamos que $m = \{1, 2, \dots\}$. Si X es un proceso de salto de Markov recurrente e irreducible, entonces independientemente de i ,

$$\lim_{t \rightarrow \infty} P(X_t = j \mid X_0 = i) = \pi_j \quad (1.41)$$

por algún número $\pi_j > 0$. Además, $\pi = (\pi_1, \pi_2, \dots)$ es la solución a

$$\sum_{j=i} \pi_j q_{j_0} = \sum_{j=i} \pi_j q_{j_1}, \text{ para todos } i = 1, \dots, \text{metro} \quad (1.42)$$

con $\sum_j \pi_j = 1$, si tal solución existe, en cuyo caso todos los estados son positivos recurrentes. Si tal solución no existe, todos π_j son 0.

Como en el caso de la cadena de Markov, $\{\pi_j\}$ se llama *limitar la distribución* $\sum_j F_j X_j$ es generalmente identificado con el vector de fila π . Alguna solución π de (1.42) con $\sum_j \pi_j = 1$ se llama *distribución estacionaria*, ya que tomarlo como la distribución inicial del proceso de salto de Markov hace que el proceso sea estacionario.

Las ecuaciones (1.42) se llaman de *nuevas ecuaciones de equilibrio global* y se generalizan fácilmente a (1.37), reemplazando las probabilidades de transición con tasas de transición. Más importante aún, si el proceso es reversible, entonces, al igual que con las cadenas de Markov, la distribución estacionaria se puede encontrar a partir de las *ecuaciones de balance local*:

$$\pi_i q_{j0} = \pi_j q_{ji}, y_0, j \in m_i \quad (1.43)$$

La reversibilidad se puede verificar fácilmente comprobando que no se produzca bucle, es decir, mediante el criterio de Kolmogorov (1.39), reemplazando las probabilidades *pag* con tarifas *q*.

■ EJEMPLO 1.13 M/M/1 cola

Considere una instalación de servicio donde los clientes llegan en ciertos momentos aleatorios y son atendidos por un solo servidor. Los clientes que llegan y encuentran el servidor ocupado esperan en la cola. Los clientes son atendidos en el orden en que llegan. Los tiempos entre llegadas son variables aleatorias exponenciales con tasas λ , y los tiempos de servicio de los clientes son variables aleatorias exponenciales iid con tasas μ . Por último, los tiempos de servicio son independientes de los tiempos entre llegadas. Dejar X_t sea el número de clientes en el sistema en ese momento t . Por la propiedad sin memoria de la distribución exponencial (vea el problema 1.7), no es difícil ver que $X = \{X_t, t \geq 0\}$ es un proceso de salto de Markov y, de hecho, un proceso de nacimiento y muerte con tasas de natalidad $b = \lambda, i = 0, 1, 2, \dots$ y tasas de mortalidad $d = \mu, i = 1, 2, \dots$.

Resolver las ecuaciones de balance global (o, más fácilmente, las ecuaciones de balance local, ya que X es reversible), vemos que X tiene una distribución límite dada por

$$\lim_{t \rightarrow \infty} \text{PAGS}(X_t = \text{norte}) = (1 - \rho)^{\text{norte}}, \text{norte} = 0, 1, 2, \dots, \quad (1.44)$$

siempre que $\rho = \lambda/\mu < 1$. Esto significa que el tiempo de servicio esperado debe ser menor que el tiempo entre llegadas esperado para que exista una distribución límite. En ese caso, la distribución límite es también la distribución estacionaria. En particular, si X_0 se distribuye de acuerdo con (1.44), entonces X_t tiene la misma distribución para todos $t > 0$.

1.14 PROCESOS GAUSSIANOS

La distribución normal también se llama *distribución gaussiana* o *distribución*. Los procesos gaussianos son generalizaciones de vectores aleatorios normales multivariados (discutidos en la Sección 1.10). Específicamente, un proceso estocástico $\{X_t, t \in T\}$ se ha dicho *gaussiano* si todas sus distribuciones de dimensión finita son gaussianas. Es decir, si para cualquier elección de $\text{norte} y t_1, \dots, t_{\text{norte}} \in T$, sostiene que

$$(X_{t_1}, \dots, X_{t_{\text{norte}}}) \sim \text{NORTE}(\mathbf{m}, \Sigma) \quad (1.45)$$

para algún vector de expectativa \mathbf{m} y matriz de covarianza Σ (ambos dependen de la elección Σ de $t_1, \dots, t_{\text{norte}}$). ~~combinación~~ Equivalentemente, $\{X_t, t \in T\}$ es gaussiana si alguna lineal $\sum_{i=1}^{\text{norte}} b_i X_{t_i}$ tiene una distribución normal. Tenga en cuenta que un proceso gaussiano está completamente determinada por su *función de expectativa* $\mu_t = \mathbb{E}[X_t]$, $t \in T$, y *función de covarianza* $\Sigma_{st} = \text{cov}(X_s, X_t)$, $s, t \in T$.

■ EJEMPLO 1.14 Proceso de Wiener (movimiento browniano)

El proceso gaussiano por excelencia es el *proceso de salchichao* (estándar) *movimiento browniano*. Puede verse como una versión continua de un proceso de caminata aleatoria. La Figura 1.8 da una ruta de muestra típica. El proceso de Wiener juega un papel central en la probabilidad y constituye la base de muchos otros procesos estocásticos.

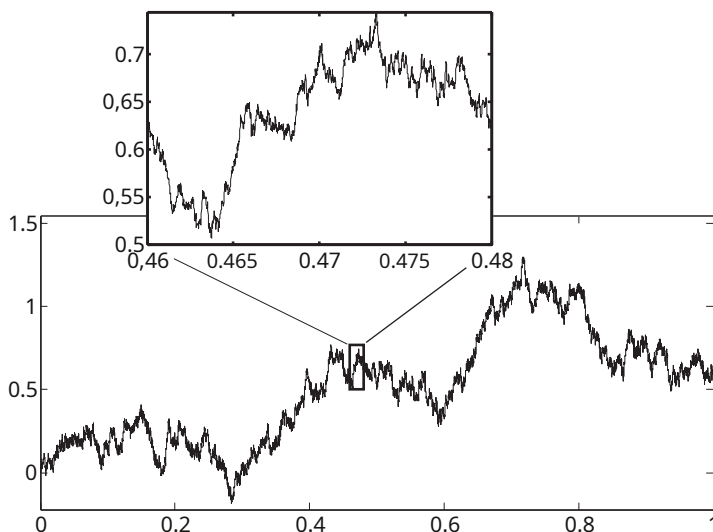


Figura 1.8: Una ruta de muestra del proceso de Wiener. El recuadro muestra una ampliación de la ruta durante un pequeño intervalo de tiempo.

El proceso de Wiener se puede definir como un proceso gaussiano $\{X_t, t \geq 0\}$ con función de expectativa $m_t = 0$ para todo t y función de covarianza $\Sigma_{s,t} = s$ para $0 \leq s \leq t$. El proceso de Wiener tiene muchas propiedades fascinantes (p. ej., [11]). Por ejemplo, es un proceso de Markov (es decir, satisface la propiedad de Markov (1.30)) con rutas de muestra continuas que son en *ninguna parte diferenciable*. Además, los incrementos $X_t - X_s$ en intervalos $[s, t]$ son independientes y normalmente distribuidos. En concreto, para cualquier $t_1 < t_2 < t_3 < t_4$,

$$X_{t_4} - X_{t_3} \quad \text{y} \quad X_{t_2} - X_{t_1}$$

son variables aleatorias independientes, y para todos $t \geq s \geq 0$,

$$X_t - X_s \sim \text{NORTE}(0, t - s).$$

Esto conduce a un procedimiento de simulación simple para los procesos de Wiener, que se analiza en la Sección 2.8.

1.15 INFORMACIÓN

En esta sección discutimos brevemente varias medidas de información en un experimento aleatorio. Supongamos que describimos las medidas en un experimento aleatorio a través de

un vector aleatorio $\mathbf{X} = (X_1, \dots, X_{norte})$ con pdf F . Entonces toda la información sobre el experimento (todo nuestro conocimiento probabilístico) obviamente está contenida en el pdf F . Sin embargo, en la mayoría de los casos nos gustaría caracterizar nuestra información sobre los experimentos con solo unos pocos números clave, como el *expectativa* y la *Matriz de covarianza* de \mathbf{X} , que proporcionan información sobre las medidas medias y la variabilidad de las medidas, respectivamente. Otra medida informativa proviene de la teoría de la codificación y las comunicaciones, donde el *Entropía de Shannon* caracteriza el número promedio de bits necesarios para transmitir un mensaje \mathbf{X} a través de un canal de comunicación (binario). Otro enfoque más de la información se puede encontrar en las estadísticas. Específicamente, en la teoría de la estimación puntual, la función de densidad de probabilidad f depende de un vector de parámetros θ . La pregunta es qué tan bien θ se puede estimar a través de un resultado de \mathbf{X} —en otras palabras, cuánta información sobre θ está contenido en los "datos" \mathbf{X} . Diversas medidas para este tipo de información están asociadas a la *máxima verosimilitud*, la *puntaje*, y el (*Fisher*) *matriz de información*. Finalmente, la cantidad de información en un experimento aleatorio a menudo se puede cuantificar a través de un *distancia* concepto, como el *Kullback-Leibler* "distancia" (divergencia), también llamada *entropía cruzada*.

1.15.1 Entropía de Shannon

Una de las medidas de incertidumbre más celebradas en la teoría de la información es la *Entropía de Shannon*, o simplemente *entropía*. Una buena referencia es [4], donde la entropía de una variable aleatoria discreta X con densidad f se define como

$$H(X) = -\sum_{x \in \mathcal{X}} f(x) \ln f(x) = -\sum_{x \in \mathcal{X}} P(X=x) \ln P(X=x).$$

Aquí X se interpreta como un carácter aleatorio de un alfabeto \mathcal{X} , tal que $X \sim \mathcal{X}$ con probabilidad $f(x)$. Usaremos la convención $0 \ln 0 = 0$.

Se puede demostrar que la forma más eficiente de transmitir caracteres muestreados de f sobre un canal binario es codificarlos de tal manera que el \sum número de bits necesarios transmitir X es igual a $\log_2(1/f(X))$. Resulta que $\sum_{x \in \mathcal{X}} f(x) \log_2(1/f(x))$ es la longitud de bit esperada necesaria para enviar un carácter aleatorio $X \sim f$; ver [4].

Un enfoque más general, que incluye variables aleatorias continuas, es definir la entropía de una variable aleatoria X con densidad f por

$$H(X) = -\mathbb{E}[\ln f(X)] = \begin{cases} -\sum_{x \in \mathcal{X}} f(x) \ln f(x) & \text{caso discreto,} \\ -\int_{\mathcal{X}} f(x) \ln f(x) dx & \text{caso continuo.} \end{cases} \quad (1.46)$$

La definición (1.46) se puede extender fácilmente a vectores aleatorios \mathbf{X} como (en el caso continuo)

$$H(\mathbf{X}) = -\mathbb{E}[\ln f(\mathbf{X})] = -\int_{\mathcal{X}} f(\mathbf{x}) \ln f(\mathbf{x}) d\mathbf{x}. \quad (1.47)$$

$H(\mathbf{X})$ es a menudo llamado el *articulación* entropía de las variables aleatorias X_1, \dots, X_{norte} , y también se escribe como $H(X_1, \dots, X_{norte})$. En el caso continuo, $H(\mathbf{X})$ es frecuentemente referido como el *entropía diferencial* distinguirlo del caso discreto.

■ EJEMPLO 1.15

Dejar X tener un $\text{Ber}(pags)$ distribución para algunos $0 < pags < 1$. La densidad F de X es dado por $F(1) = \text{PAGS}(X=1) = pags$ y $F(0) = \text{PAGS}(X=0) = 1 - pags$ de modo que la entropía de X es

$$H(X) = -pags \ln pags - (1 - pags) \ln (1 - pags).$$

La gráfica de la entropía en función de $pags$ se representa en la figura 1.9. Tenga en cuenta que la entropía es máxima para $pags = 1/2$, que da la densidad "uniforme" en $\{0, 1\}$.

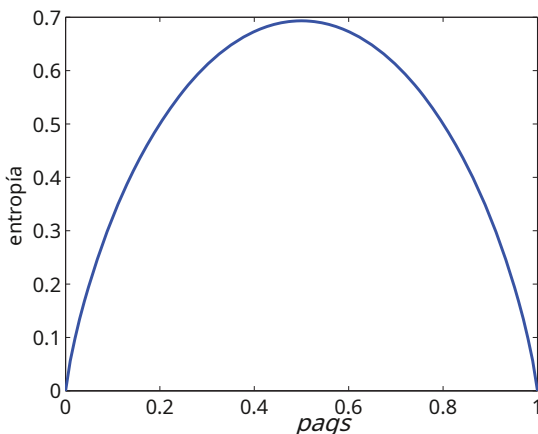


Figura 1.9: La entropía para el $\text{Ber}(pags)$ distribución en función de $pags$.

A continuación, considere una secuencia X_1, \dots, X_{norte} de iid $\text{Ber}(pags)$ variables aleatorias. Dejar $\mathbf{X} = (X_1, \dots, X_{norte})$. La densidad de \mathbf{X} , decir *gramo*, es simplemente el producto de las densidades de los X_i , de modo que

$$H(\mathbf{X}) = -\text{mi} [\text{engramo}(\mathbf{X})] = -\text{mi} \left[\prod_{i=1}^{norte} F(X_i) \right] = \sum_{i=1}^{norte} -\text{mi} [\ln F(X_i)] = norte H(X).$$

las propiedades de $H(\mathbf{X})$ en el caso continuo son algo diferentes de las del caso discreto. En particular:

1. La entropía diferencial puede ser negativa, mientras que la entropía discreta siempre es positiva.
2. La entropía discreta es insensible a las transformaciones invertibles, mientras que la entropía diferencial no lo es. Específicamente, si X es discreto, $Y = \text{gramo}(X)$, y gramo es una aplicación invertible, entonces $H(X) = H(Y)$ porque $F_Y(y) = F_X(\text{gramo}^{-1}(y))$. Sin embargo, en el caso continuo, tenemos un factor adicional debido al jacobiano de la transformación.

No es difícil ver que de cualquier densidad F , la que da la máxima entropía es la densidad uniforme en X . Eso es,

$$H(\mathbf{X}) \text{ es máximo} \Leftrightarrow F(\mathbf{X}) = \frac{1}{|J_X|} \text{ (constante)}. \quad (1.48)$$

Para dos vectores aleatorios \mathbf{X} y \mathbf{Y} con pdf conjunto F , definimos la *entropía condicional de \mathbf{Y} dado \mathbf{X}* como

$$H(\mathbf{Y}/\mathbf{X}) = -\text{mien} \left[\frac{F(\mathbf{X}, \mathbf{Y})}{F_{\mathbf{X}}(\mathbf{X})} \right] = H(\mathbf{X}, \mathbf{Y}) - H(\mathbf{X}), \quad (1.49)$$

dónde $F_{\mathbf{X}}$ es el pdf de \mathbf{X} y $F(\mathbf{X}, \mathbf{y})$ es la densidad condicional de \mathbf{Y} (ay), dado $\mathbf{X}=\mathbf{X}$. Resulta que

$$H(\mathbf{X}, \mathbf{Y}) = H(\mathbf{X}) + H(\mathbf{Y}/\mathbf{X}) = H(\mathbf{Y}) + H(\mathbf{X}/\mathbf{Y}). \quad (1.50)$$

Es razonable exigir que cualquier medida aditiva sensata que describa la cantidad promedio de incertidumbre satisfaga al menos (1.50) y (1.48). De ello se deduce que la densidad uniforme lleva la menor cantidad de información, y la entropía (cantidad media de incertidumbre) de (\mathbf{X}, \mathbf{Y}) es igual a la suma de la entropía de \mathbf{X} y la cantidad de entropía en \mathbf{Y} después de la información en \mathbf{X} ha sido contabilizado. Se argumenta en [10] que cualquier concepto de entropía que incluya las propiedades generales (1.48) y (1.50) debe conducir a la definición (1.47).

la *información mutua* de \mathbf{X} y \mathbf{Y} Se define como

$$\text{METRO}(\mathbf{X}, \mathbf{Y}) = H(\mathbf{X}) + H(\mathbf{Y}) - H(\mathbf{X}, \mathbf{Y}), \quad (1.51)$$

que, como su nombre indica, puede interpretarse como la cantidad de información compartida por \mathbf{X} y \mathbf{Y} . Una expresión alternativa, que se sigue de (1.50) y (1.51), es

$$\text{METRO}(\mathbf{X}, \mathbf{Y}) = H(\mathbf{X}) - H(\mathbf{X}/\mathbf{Y}) = H(\mathbf{Y}) - H(\mathbf{Y}/\mathbf{X}), \quad (1.52)$$

que puede interpretarse como la reducción de la incertidumbre de una variable aleatoria debido al conocimiento de la otra. No es difícil demostrar que la información mutua siempre es positiva. También está relacionado con el concepto de entropía cruzada, que sigue.

1.15.2 Entropía cruzada de Kullback-Leibler

Dejar g y h ser dos densidades en \mathbf{X} . La entropía cruzada de Kullback-Leibler entre g y h (comparar con (1.47)) se define (en el caso continuo) como

$$\begin{aligned} D(g, h) &= \int \left[\frac{g(\mathbf{X})}{h(\mathbf{X})} \right] \ln \left[\frac{g(\mathbf{X})}{h(\mathbf{X})} \right] h(\mathbf{X}) d\mathbf{X} \\ &= \int g(\mathbf{X}) \ln \left[\frac{g(\mathbf{X})}{h(\mathbf{X})} \right] d\mathbf{X} = \int g(\mathbf{X}) \ln g(\mathbf{X}) d\mathbf{X} - \int g(\mathbf{X}) \ln h(\mathbf{X}) d\mathbf{X}. \end{aligned} \quad (1.53)$$

$D(g, h)$ también se llama *Divergencia Kullback-Leibler*, la *entropía cruzada*, y el *entropía relativa*. Si no se indica lo contrario, llamaremos $D(g, h)$ la *entropía cruzada* (CE) entre g y h . Darse cuenta de $D(g, h)$ no es una distancia entre g y h en el sentido formal, ya que en general $D(g, h) \neq D(h, g)$. Sin embargo, a menudo es útil pensar en $D(g, h)$ como una distancia porque

$$D(g, h) \geq 0$$

y $D(g, h) = 0$ si y solo si $\text{gramo}(X) = h(X)$. Esto se sigue de la desigualdad de Jensen (si φ es una función convexa, tal como $-\ln$ entonces $MI[\varphi(X)] \geq \varphi(MI[X])$). A saber

$$D(g, h) = MI_{\text{gramo}(X)} \left[\frac{h(X)}{\text{gramo}(X)} \right] - \ln \left(MI_{\text{gramo}(X)} \left[\frac{h(X)}{\text{gramo}(X)} \right] \right) = -\ln 1 = 0.$$

Se puede ver fácilmente que la información mutua $METRO(X, Y)$ de vectores X y Y definido en (1.51) está relacionado con el CE de la siguiente manera:

$$METRO(X, Y) = D(f_{XY}, f_X f_Y) = MI_{f_X(X) f_Y(Y)} \left[\frac{f_{XY}(X, Y)}{f_X(X) f_Y(Y)} \right],$$

dónde f es el pdf (conjunto) de (X, Y) y f_X, f_Y son los pdf (marginales) de X y Y , respectivamente. En otras palabras, la información mutua se puede ver como el CE que mide la distancia entre el pdf conjunto f_{XY} y el producto de sus pdf marginales $f_X f_Y$, es decir, bajo el supuesto de que los vectores X y Y son *independiente*.

1.15.3 Estimador de máxima verosimilitud y función de puntuación

Introducimos aquí la noción de *función de puntuación* (SF) a través de la clásica *estimador de máxima verosimilitud*. Considere un vector aleatorio $X = (X_1, \dots, X_{norte})$ que se distribuye de acuerdo con un pdf fijo $f; \theta$ con parámetro desconocido (vector) $\theta \in \Theta$. Digamos que queremos estimar θ sobre la base de un resultado dado X (los datos) de X . Para una dada X , la función $L(\theta; X) = f(X; \theta)$ se llama *función de probabilidad*. Tenga en cuenta que L es una función de θ para un parámetro fijo X , mientras que para el pdf f es al revés. La máxima probabilidad *estimar* $\hat{\theta} = \hat{\theta}(X)$ de θ se define como

$$\hat{\theta} = \arg\max_{\theta \in \Theta} L(\theta; X). \quad (1.54)$$

Como la función \ln es monótona creciente, también tenemos

$$\hat{\theta} = \arg\max_{\theta \in \Theta} \ln L(\theta; X). \quad (1.55)$$

la variable aleatoria $\hat{\theta}(X)$ con $X \sim f; \theta$ es la probabilidad máxima correspondiente *capucha estimador*, que se escribe de nuevo como $\hat{\theta}$. Tenga en cuenta que a menudo el $\hat{\theta}$ ejército de reserva X_1, \dots, X_{norte} formar una muestra aleatoria de algún pdf $f; \theta$, en ese caso $f(X; \theta) = \prod_{i=1}^{norte} f(X_i; \theta)$

$$\hat{\theta} = \arg\max_{\theta \in \Theta} \sum_{i=1}^{norte} \ln f(X_i; \theta). \quad (1.56)$$

Si $L(\theta; X)$ es una función cóncava continuamente diferenciable con respecto a θ y el máximo se alcanza en el interior de Θ , entonces podemos encontrar el estimador de máxima verosimilitud de θ resolviendo

$$\nabla_{\theta} \ln L(\theta; X) = 0.$$

La función $S(\cdot; X)$ definido por

$$S(\theta; X) = \nabla_{\theta} \ln L(\theta; X) = \frac{\nabla_{\theta} f(X; \theta)}{f(X; \theta)} \quad (1.57)$$

se llama *función de puntuación*. Para la familia exponencial (A.9) es fácil ver que

$$S(\boldsymbol{\theta}; \mathbf{X}) = \frac{\nabla \zeta(\boldsymbol{\theta})}{\zeta(\boldsymbol{\theta})} + \mathbf{t}(\mathbf{X}). \quad (1.58)$$

El *vector aleatorio* $S(\boldsymbol{\theta}) = S(\boldsymbol{\theta}; \mathbf{X})$ con $\mathbf{X} \sim \mathcal{P}(\cdot; \boldsymbol{\theta})$ se llama el *eficiente* *puntuación*. La puntuación esperada siempre es igual al vector cero, es decir,

$$\mathbb{E}_{\boldsymbol{\theta}}[S(\boldsymbol{\theta})] = \int \nabla_{\boldsymbol{\theta}} \mathcal{P}(\mathbf{X}; \boldsymbol{\theta}) m(d\mathbf{X}) = \nabla_{\boldsymbol{\theta}} \int \mathcal{P}(\mathbf{X}; \boldsymbol{\theta}) m(d\mathbf{X}) = \nabla_{\boldsymbol{\theta}} 1 = \mathbf{0},$$

donde el intercambio de diferenciación e integración se justifica a través del teorema de convergencia acotada.

1.15.4 Información del pescador

La matriz de covarianza $\text{YO}(\boldsymbol{\theta})$ de la partitura $S(\boldsymbol{\theta})$ se llama *Matriz de información de Fisher*. Dado que la puntuación esperada siempre es $\mathbf{0}$, tenemos

$$\text{YO}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}}[S(\boldsymbol{\theta}) S(\boldsymbol{\theta})^T]. \quad (1.59)$$

En el caso unidimensional, tenemos entonces

$$\text{YO}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}} \left[\left(\frac{\partial \ln \mathcal{P}(\mathbf{X}; \boldsymbol{\theta})}{\partial \theta} \right)^2 \right].$$

Porque

$$\frac{\partial^2}{\partial \theta^2} \ln \mathcal{P}(\mathbf{X}; \boldsymbol{\theta}) = \frac{\frac{\partial^2}{\partial \theta^2} \mathcal{P}(\mathbf{X}; \boldsymbol{\theta})}{\mathcal{P}(\mathbf{X}; \boldsymbol{\theta})} - \left(\frac{\frac{\partial}{\partial \theta} \mathcal{P}(\mathbf{X}; \boldsymbol{\theta})}{\mathcal{P}(\mathbf{X}; \boldsymbol{\theta})} \right)^2,$$

vemos que (en condiciones de regularidad directa) la información de Fisher también está dada por

$$\text{YO}(\boldsymbol{\theta}) = -\mathbb{E}_{\boldsymbol{\theta}} \left[\frac{\partial^2 \ln \mathcal{P}(\mathbf{X}; \boldsymbol{\theta})}{\partial \theta^2} \right].$$

En el caso multidimensional tenemos de manera similar

$$\text{YO}(\boldsymbol{\theta}) = -\mathbb{E}_{\boldsymbol{\theta}}[\nabla S(\boldsymbol{\theta})] = -\mathbb{E}_{\boldsymbol{\theta}} \left[\nabla_{\boldsymbol{\theta}} \ln \mathcal{P}(\mathbf{X}; \boldsymbol{\theta}) \right] \quad (1.60)$$

dónde $\nabla_{\boldsymbol{\theta}} \ln \mathcal{P}(\mathbf{X}; \boldsymbol{\theta})$ denota el *arpillerado* $\ln \mathcal{P}(\mathbf{X}; \boldsymbol{\theta})$, es decir, la matriz (aleatoria)

$$\left(\frac{\partial^2 \ln \mathcal{P}(\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right).$$

La importancia de la información de Fisher en las estadísticas es corroborada por el famoso *Desigualdad de Cramér-Rao*, que (de forma simplificada) establece que la varianza de cualquier estimador insesgado Z de $\text{gramo}(\boldsymbol{\theta})$ está acotado desde abajo a través de

$$\text{Var}(Z) \geq (\nabla \text{gramo}(\boldsymbol{\theta}))_{\text{YO}^{-1}(\boldsymbol{\theta})} \nabla \text{gramo}(\boldsymbol{\theta}). \quad (1.61)$$

Para más detalles, véase [12].

1.16 OPTIMIZACIÓN CONVEXA Y DUALIDAD

Dejar $f(X)$, $X \in \mathbb{R}^n$, sea una función de valor real con derivadas continuas, también llamada función C^1 función. El enfoque estándar para minimizar $f(X)$ es resolver la ecuación

$$\nabla f(X) = 0. \quad (1.62)$$

Las soluciones de (1.62) se llaman *puntos estacionarios*. Si, además, la función tiene segundas derivadas continuas (la llamada C^2 función), la condición

$$\nabla^2 f(X^*) > 0 \quad (1.63)$$

asegura que un punto estacionario X^* es un *minimizador local*, eso es, $f(X^*) < f(X)$ para todos X en un barrio bastante pequeño de X^* .

Para C^1 función en \mathbb{R}^n , (1.62) se generaliza a

$$\nabla f(X) = \begin{pmatrix} \frac{\partial f(X)}{\partial x_1} \\ \vdots \\ \frac{\partial f(X)}{\partial x_n} \end{pmatrix} = 0, \quad (1.64)$$

dónde $\nabla f(X)$ es el *gradiente* de f en X . Del mismo modo, un punto estacionario X^* es un *minimizador local* de f si la *matriz Hessiana* (o simplemente *arpillera*) en X^* ,

$$\nabla^2 f(X^*) = \begin{pmatrix} \frac{\partial^2 f(X^*)}{\partial x_1^2} & \dots & \frac{\partial^2 f(X^*)}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(X^*)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(X^*)}{\partial x_n^2} \end{pmatrix}, \quad (1.65)$$

es *positivo definitivo*, eso es, $X^T [\nabla^2 f(X^*)] X > 0$ para todos $X \neq 0$.

La situación se puede generalizar aún más introduciendo *restricciones*. Un problema general de optimización con restricciones se puede escribir como

$$\min_{X \in \mathbb{R}^n} f(X) \quad (1.66)$$

$$\text{sujeto a: } h_i(X) = 0, \quad i = 1, \dots, m, \quad (1.67)$$

$$g_j(X) \leq 0, \quad j = 1, \dots, k \quad (1.68)$$

Aquí, f , g_i , y h_j son funciones, $f(X)$ se llama la *función objetivo*, y $h_i(X) = 0$ y $g_j(X) \leq 0$ representan las *igualdades* y *desigualdades* restricciones, respectivamente.

La región del dominio donde está definida la función objetivo y donde se satisfacen todas las restricciones se llama *región factible*. Una *solución mundial* al problema de optimización es un punto $X^* \in \mathbb{R}^n$ tal que no existe otro punto $X \in \mathbb{R}^n$ para cual $f(X) < f(X^*)$. Los nombres alternativos son *minimizador global* y *mínimo global*, aunque este último podría confundirse con el valor mínimo de la función. Del mismo modo, para un *local* solución/minimizador, la condición $f(X) < f(X^*)$ solo necesita mantenerse en algún vecindario de X^* .

Dentro de esta formulación se encuentran muchos de los problemas de optimización tradicionales. Un problema de optimización en el que la función objetivo y las restricciones de igualdad y desigualdad son funciones lineales se denomina *programa lineal*. Una optimización

problema en el que la función objetivo es cuadrática, mientras que las restricciones son funciones lineales se llama un *programa cuadrático*. La convexidad juega un papel importante en muchos problemas prácticos de optimización.

Definición 1.16.1 (Conjunto convexo) Un conjunto $X \in \mathbb{R}^n$ se llama *convexo* si por todo $\mathbf{x}, \mathbf{y} \in X$ $\theta \in (0, 1)$, el punto $(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \in X$.

Definición 1.16.2 (Función convexa) Una función $f(\mathbf{x})$ en un conjunto convexo X se llama *convexo* si por todo $\mathbf{x}, \mathbf{y} \in X$ $\theta \in (0, 1)$,

$$f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y}). \quad (1.69)$$

Si se cumple una desigualdad estricta en (1.69), se dice que la función es *estrictamente convexo*. Si una función f es (estrictamente) convexo, entonces $-f$ es (estrictamente) *cóncavo*. Asumiendo X es un conjunto abierto, convexidad para $f \in C_1$ es equivalente a

$$f(\mathbf{y}) - f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) \geq 0 \quad \text{para todos } \mathbf{x}, \mathbf{y} \in X.$$

Además, por $f \in C_2$, la convexidad es equivalente a que la matriz hessiana sea semidefinida positiva para todo $\mathbf{x} \in X$:

$$\mathbf{y}^T [\nabla^2 f(\mathbf{x})] \mathbf{y} \geq 0, \text{ para todos } \mathbf{y} \in \mathbb{R}^n.$$

Se dice que el problema (1.66) es un *problema de programación convexa* si

1. la función objetivo f es convexo,
2. las funciones de restricción de desigualdad $\{g_i(\mathbf{x})\}$ son convexas y
3. las funciones de restricción de igualdad $\{h_i(\mathbf{x})\}$ son *afín*, es decir, de la forma $\mathbf{a}_i^T \mathbf{x} - b_i$.

Tenga en cuenta que el último requisito se deriva del hecho de que una restricción de igualdad $h_i(\mathbf{x}) = 0$ puede verse como una combinación de las restricciones de desigualdad $h_i(\mathbf{x}) - 0$ y $-h_i(\mathbf{x}) - 0$, de modo que ambos h_i necesita ser convexo. Tanto el programa lineal como el cuadrático (con matriz definida positiva C) son convexas.

1.16.1 Método Lagrangiano

Los componentes principales del método de Lagrange son los multiplicadores de Lagrange y la función de Lagrange. El método fue desarrollado por Lagrange en 1797 para el problema de optimización (1.66) con restricciones de igualdad (1.67). En 1951 Kuhn y Tucker extendió el método de Lagrange a las restricciones de desigualdad.

Definición 1.16.3 (Función de Lagrange) Dado un problema de optimización (1.66) que contiene solo restricciones de igualdad $h_i(\mathbf{x}) = 0, i = 1, \dots, m$, la *Función de Lagrange*, o *Lagrangiano*, Se define como

$$L(\mathbf{x}, \boldsymbol{\beta}) = f(\mathbf{x}) + \sum_i \beta_i h_i(\mathbf{x}),$$

donde los coeficientes $\{\beta_i\}$ se llaman los *Multiplicadores de Lagrange*.

Una condición necesaria para un punto \mathbf{X}^* ser un minimizador local de $F(\mathbf{X})$ sujeto a las restricciones de igualdad $h_i(\mathbf{X}) = 0, i=1, \dots, m$, es

$$\nabla_{\mathbf{X}} L(\mathbf{X}^*, \boldsymbol{\beta}^*) = \mathbf{0}, \quad \nabla_{\boldsymbol{\beta}} L(\mathbf{X}^*, \boldsymbol{\beta}^*) = \mathbf{0},$$

por algún valor $\boldsymbol{\beta}^*$. Las condiciones anteriores también son suficientes si $L(\mathbf{X}, \boldsymbol{\beta}^*)$ es una función convexa de \mathbf{X} .

■ EJEMPLO 1.16 Distribución máxima de entropía

Dejar $\mathbf{pags} = \{pags_i, i=1, \dots, n\}$ Sea una distribución de probabilidad. Considere el siguiente programa, que maximiza la entropía (Shannon):

$$\begin{aligned} & \text{máximo} - \sum_{i=1}^{\text{norte}} pags_i \ln pags_i \\ & \text{sujeto a:} \quad \sum_{i=1}^{\text{norte}} pags_i = 1. \end{aligned}$$

El lagrangiano es

$$L(\mathbf{pags}, \beta) = \sum_{i=1}^{\text{norte}} pags_i \ln pags_i + \beta \left(\sum_{i=1}^{\text{norte}} pags_i - 1 \right)$$

sobre el dominio $\{(\mathbf{pags}, \beta) : pags_i \geq 0, i=1, \dots, \text{norte}, \beta \in \mathbb{R}\}$. La solución óptima \mathbf{pags}^* del problema es la distribución uniforme, es decir, $\mathbf{pags}^* = (1/n, \dots, 1/\text{norte})$; vea el problema 1.35.

Definición 1.16.4 (Función de Lagrange generalizada) Dado el problema de optimización original (1.66), que contiene las restricciones de igualdad y desigualdad, el *función de Lagrange generalizada*, o simplemente *Lagrangiano*, Se define como

$$L(\mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = F(\mathbf{X}) + \sum_{i=1}^k \alpha_i g_i(\mathbf{X}) + \sum_{j=1}^m \beta_j h_j(\mathbf{X}).$$

Una condición necesaria para un punto \mathbf{X}^* ser un minimizador local de $F(\mathbf{X})$ en el problema de optimización (1.66) es la existencia de un $\boldsymbol{\alpha}^*$ y $\boldsymbol{\beta}^*$ tal que

$$\begin{aligned} & \nabla_{\mathbf{X}} L(\mathbf{X}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) = \mathbf{0}, \quad \nabla_{\boldsymbol{\beta}} L(\mathbf{X}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) = \mathbf{0}, \\ & \alpha_i g_i(\mathbf{X}^*) = 0, \quad i=1, \dots, k, \\ & \alpha_i \neq 0, \quad i=1, \dots, k, \\ & \alpha_i g_i(\mathbf{X}^*) = 0, \quad i=1, \dots, k \end{aligned}$$

Estas ecuaciones generalmente se denominan *Condiciones de Karush-Kuhn-Tucker (KKT)*. Para *convexos* programas tenemos los siguientes resultados importantes:

1. Cada solución local \mathbf{X}^* a un problema de programación convexo es una solución global y el conjunto de soluciones globales es convexo. Si, además, la función objetivo es estrictamente convexa, entonces cualquier solución global es única.

2. Para un problema de programación estrictamente convexa con C funciones objetivo y de restricción, las condiciones KKT son necesarias y suficientes para una solución global única.

1.16.2 Dualidad

El objetivo de la dualidad es proporcionar una formulación alternativa de un problema de optimización que a menudo es computacionalmente más eficiente o tiene algún significado teórico (ver [7], página 219). El problema original (1.66) se conoce como el *primitivo* problema, mientras que el problema reformulado, basado en los multiplicadores de Lagrange, se conoce como el *doble* problema. La teoría de la dualidad es más relevante para los problemas de optimización convexa. Es bien sabido que si el problema de optimización primal es (estrictamente) convexo, entonces el problema dual es (estrictamente) cóncavo y tiene una solución (única) a partir de la cual se puede deducir la solución primal óptima (única).

Definición 1.16.5 (Programa dual de Lagrange) El programa dual de Lagrange del programa primal (1.66), es

$$\begin{aligned} & \underset{\alpha, \beta}{\text{máximo}} && L(\alpha, \beta) \\ \text{sujeto a:} &&& \alpha \geq 0, \end{aligned}$$

dónde L es el *doble función de Lagrange*:

$$L(\alpha, \beta) = \inf_{x \in X} L(x, \alpha, \beta). \quad (1.70)$$

No es difícil ver que si F^* es el valor mínimo del problema primal, entonces $L(\alpha, \beta) \leq F^*$ para cualquier $\alpha \geq 0$ y cualquier β . Esta propiedad se llama *dualidad débil*. El programa dual lagrangiano determina así el mejor límite inferior de F^* . Si d^* es el valor óptimo para el problema dual, entonces $d^* \leq F^*$. La diferencia $F^* - d^*$ se llama el *brecha de dualidad*.

La brecha de dualidad es extremadamente útil para proporcionar límites inferiores para las soluciones de problemas primarios que pueden ser imposibles de resolver directamente. Es importante tener en cuenta que para los problemas con restricciones lineales, si el primario no es factible (no tiene una solución que satisfaga las restricciones), entonces el dual es no factible o ilimitado. Por el contrario, si el dual es inviable, entonces el primal no tiene solución. De crucial importancia es la *fuerte dualidad* teorema, que establece que para programas convexos (1.66) con funciones lineales restringidas *hygram* la brecha de dualidad es cero, y cualquier x^* y (α^*, β^*) que satisfacen las condiciones KKT son soluciones (globales) para los programas primal y dual, respectivamente. En particular, esto es válido para los programas cuadráticos lineales y convexos (tenga en cuenta que no todos los programas cuadráticos son convexos).

Para un programa primal convexo con C funciones objetivo y de restricción, la función dual lagrangiana (1.70) se puede obtener simplemente estableciendo el gradiente (con respecto a x) del lagrangiano $L(x, \alpha, \beta)$ a cero. Se puede simplificar aún más el programa dual sustituyendo en el Lagrangiano las relaciones entre las variables así obtenidas.

■ EJEMPLO 1.17 Problema de Programación Lineal

Considere el siguiente problema de programación lineal:

$$\begin{aligned} \min_x \quad & \mathbf{c}\mathbf{x} \\ \text{sueto a:} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned}$$

El lagrangiano es $L(\mathbf{x}, \boldsymbol{\alpha}) = \mathbf{c}\mathbf{x} - \boldsymbol{\alpha}(\mathbf{A}\mathbf{x} - \mathbf{b})$. La función dual de Lagrange es el mínimo de L general \mathbf{x} ; de este modo

$$L_*(\boldsymbol{\alpha}) = \begin{cases} \mathbf{b}\boldsymbol{\alpha} & \text{si } \mathbf{A}\boldsymbol{\alpha} = \mathbf{c}, \\ -\infty & \text{de lo contrario,} \end{cases}$$

para que el programa dual de Lagrange se convierta en

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \mathbf{b}\boldsymbol{\alpha} \\ \text{sueto a:} \quad & \mathbf{A}\boldsymbol{\alpha} = \mathbf{c}, \\ & \boldsymbol{\alpha} \geq \mathbf{0}. \end{aligned}$$

Un hecho interesante a tener en cuenta aquí es que para el problema de programación lineal, el dual del problema dual siempre devuelve el problema primal.

■ EJEMPLO 1.18 Problema de programación cuadrática

Considere el siguiente problema de programación cuadrática:

$$\begin{aligned} \min_x \quad & \frac{1}{2} \mathbf{x}^T \mathbf{C} \mathbf{x} \\ \text{sueto a:} \quad & \mathbf{C} \mathbf{x} = \mathbf{b}, \end{aligned}$$

donde el $\mathbf{n} \times \mathbf{n}$ *n*ortematriz \mathbf{C} se supone que es definida positiva (para un problema general de programación cuadrática, la matriz \mathbf{C} siempre se puede suponer que es simétrica, pero no es necesariamente definida positiva). El lagrangiano es $L(\mathbf{x}, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{x}^T \mathbf{C} \mathbf{x} - \boldsymbol{\alpha}^T (\mathbf{C} \mathbf{x} - \mathbf{b})$. Podemos minimizar esto tomando su gradiente con respecto a \mathbf{x} y ponerlo a cero. Esto da $\mathbf{C} \mathbf{x} - \mathbf{C} \boldsymbol{\alpha} = \mathbf{C}(\mathbf{x} - \boldsymbol{\alpha}) = \mathbf{0}$. La certeza positiva de \mathbf{C} implica que $\mathbf{x} = \boldsymbol{\alpha}$. La maximización de la Lagrangiano ahora se reduce a maximizar $L(\boldsymbol{\alpha}, \boldsymbol{\alpha}) = \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{C} \boldsymbol{\alpha} - \boldsymbol{\alpha}^T (\mathbf{C} \boldsymbol{\alpha} - \mathbf{b}) = -\frac{1}{2} \boldsymbol{\alpha}^T \mathbf{C} \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{b}$ sueto a $\boldsymbol{\alpha} \geq \mathbf{0}$. Por lo tanto, podemos escribir el problema dual como

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & -\frac{1}{2} \boldsymbol{\alpha}^T \mathbf{C} \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{b} \\ \text{sueto a:} \quad & \boldsymbol{\alpha} \geq \mathbf{0}. \end{aligned}$$

Observe que el problema dual involucra solo restricciones simples de no negatividad.

Ahora supongamos que nos dan la factorización de Cholesky $\mathbf{C} = \mathbf{L}\mathbf{L}^T$ *cam*a y *desayuno*. Resulta (vea el problema 1.36) que el dual de Lagrange del problema dual anterior se puede escribir como

$$\begin{aligned} \min_m \quad & \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\mu} \\ \text{sueto a:} \quad & \mathbf{B} \boldsymbol{\mu} = \mathbf{b}, \end{aligned} \tag{1.71}$$

con $m = B\alpha$. Este es un llamado *menor distancia* problema, el cual, siempre que conozcamos la factorización de Cholesky de C , es más fácil de resolver que el problema original de programación cuadrática.

Un último ejemplo de dualidad lo proporciona el ampliamente utilizado *método de mínima entropía cruzada* [9].

■ EJEMPLO 1.19 Método de entropía cruzada mínima (MinxEnt)

Dejar \mathbf{X} ser una variable aleatoria discreta (o vector) tomando valores $\mathbf{X}_1, \dots, \mathbf{X}_r$, y deja $\mathbf{q} = (q_1, \dots, q_r)$ y $\mathbf{pags} = (pags_1, \dots, pags_r)$ sean dos vectores de distribución (columna) estrictamente positivos para \mathbf{X} . Considere la función $f(\mathbf{pags}, \mathbf{q}) = \sum_{k=1}^r pags_k \ln(pags_k/q_k)$ siguiendo el programa primario de minimizando la entropía cruzada de \mathbf{pags} y \mathbf{q} , eso es, $\min_{\mathbf{pags}} \sum_{k=1}^r pags_k \ln(pags_k/q_k)$, para un fijo \mathbf{q} , sujeto a restricciones de igualdad lineal:

$$\min_{\mathbf{pags}} \sum_{k=1}^r pags_k \ln \frac{pags_k}{q_k} \quad (1.72)$$

$$\text{sujeto a: } \min_{\mathbf{pags}} [S(\mathbf{X})] = \sum_{k=1}^r S(\mathbf{X}_k) pags_k = y_i, \quad i=1, \dots, metro \quad (1.73)$$

$$\sum_{k=1}^r pags_k = 1, \quad (1.74)$$

dónde S_1, \dots, S_{metro} son funciones arbitrarias.

Aquí la función objetivo es convexa, ya que es una combinación lineal de funciones de la forma $pags \ln(pags)$ (ordenador personal), que son convexas en \mathbb{R}_+ , para cualquier $do > 0$. Además, las funciones de restricción de igualdad son afines (de la forma $\mathbf{a}^T \mathbf{pags} - y$). Por lo tanto, este problema es convexo. Para obtener la solución óptima \mathbf{pags}^* del programa primario anterior, normalmente es más fácil resolver el problema asociado *doble* programa [9]. A continuación presentamos el procedimiento correspondiente.

1. El Lagrangiano del problema primal viene dado por

$$L(\mathbf{pags}, \boldsymbol{\lambda}, \beta) = \sum_{k=1}^r pags_k \ln \frac{pags_k}{q_k} - \sum_{i=1}^{metro} \lambda_i \left(\sum_{k=1}^r S(\mathbf{X}_k) pags_k - y_i \right) + \beta \left(\sum_{k=1}^r pags_k - 1 \right) \quad (1.75)$$

dónde $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_{metro})$ es el vector multiplicador de Lagrange correspondiente a (1.73) y β es el multiplicador de Lagrange correspondiente a (1.74). Tenga en cuenta que podemos usar un signo más o menos en la segunda suma de (1.75). Elegimos este último porque luego generalizamos el mismo problema a la desigualdad. (-) restricciones en (1.73), dando lugar a un signo menos en el Lagrangiano.

2. Resolver (para fijo $\boldsymbol{\lambda}, \beta$)

$$\min_{\mathbf{pags}} L(\mathbf{pags}, \boldsymbol{\lambda}, \beta) \quad (1.76)$$

resolviendo

$$\nabla_{\mathbf{pags}} L(\mathbf{pags}, \boldsymbol{\lambda}, \beta) = \mathbf{0},$$

que da el conjunto de ecuaciones

$$\nabla_{\mathbf{pags}} L(\mathbf{pags}, \boldsymbol{\lambda}, \beta) = \ln \frac{pags_k}{q_k} - \sum_{i=1}^{metro} \lambda_i S(\mathbf{X}_k) + \beta = 0, \quad k=1, \dots, r.$$

Denote la solución óptima y el valor de la función óptima obtenidos del programa (1.76) como $\mathbf{pags}(\lambda, \beta)$ y $L(\lambda, \beta)$, respectivamente. Esta última es la función dual de Lagrange. Entonces escribimos

$$\mathbf{pags}_k(\lambda, \beta) = q_k \exp\{-\beta - 1 + \sum_{j=1}^{metro} \lambda_j S(\mathbf{X}_k)\}, k=1, \dots, r. \quad (1.77)$$

Dado que la suma de los $\{\mathbf{pags}_k\}$ debe ser 1, obtenemos

$$\sum_{k=1}^r q_k \exp\{-\beta - 1 + \sum_{j=1}^{metro} \lambda_j S(\mathbf{X}_k)\} = 1. \quad (1.78)$$

Sustituyendo $\mathbf{pags}(\lambda, \beta)$ de vuelta al Lagrangiano da

$$L(\lambda, \beta) = -1 + \sum_{j=1}^{metro} \lambda_j y_j - \beta. \quad (1.79)$$

3. Resuelve el *doble* programa

$$\max_{\lambda, \beta} L(\lambda, \beta). \quad (1.80)$$

Ya que β y λ están relacionados a través de (1.78), podemos resolver (1.80) sustituyendo el correspondiente $\beta(\lambda)$ en (1.79) y optimizando la función resultante:

$$D(\lambda) = -1 + \sum_{j=1}^{metro} \lambda_j y_j - \ln \left\{ \sum_{k=1}^r q_k \exp\{-1 + \sum_{j=1}^{metro} \lambda_j S(\mathbf{X}_k)\} \right\}. \quad (1.81)$$

Ya que $D(\lambda)$ es continuamente diferenciable y cóncava con respecto a λ , podemos derivar la solución óptima, λ^* , al resolver

$$\nabla_{\lambda} D(\lambda) = \mathbf{0}, \quad (1.82)$$

que se puede escribir por componentes en la siguiente forma explícita:

$$\begin{aligned} \nabla_{\lambda} D(\lambda) = y_j - & \frac{\sum_{k=1}^r S(\mathbf{X}_k) q_k \exp\{-1 + \sum_{j=1}^{metro} \lambda_j S(\mathbf{X}_k)\}}{\sum_{k=1}^r q_k \exp\{-1 + \sum_{j=1}^{metro} \lambda_j S(\mathbf{X}_k)\}} \\ = y_j - & \frac{[\sum_{k=1}^r S(\mathbf{X}_k) \exp\{-1 + \sum_{j=1}^{metro} \lambda_j S(\mathbf{X}_k)\}]}{[\sum_{k=1}^r \exp\{-1 + \sum_{j=1}^{metro} \lambda_j S(\mathbf{X}_k)\}]} = 0 \end{aligned} \quad (1.83)$$

por $j=1, \dots, metro$. El vector óptimo $\lambda^* = (\lambda_1^*, \dots, \lambda_{metro}^*)$ puede ser encontrado por resolviendo (1.83) numéricamente. Tenga en cuenta que si el programa primal tiene una solución óptima interior no vacía, entonces el programa dual tiene una solución óptima λ^* .

4. Finalmente, sustituya $\lambda = \lambda^*$ y $\beta = \beta(\lambda^*)$ de nuevo en (1.77) para obtener la solución del programa MinxEnt original.

Es importante señalar que no es necesario imponer explícitamente las condiciones $\mathbf{pags}_k \geq 0, k=1, \dots, r$, porque las cantidades $\{\mathbf{pags}_k\}$ en (1.77) son automáticamente estrictamente positivos. Esta es una propiedad crucial de la distancia CE; ver también

[1]. Es instructivo (vea el problema 1.37) verificar cómo la adición de las restricciones de no negatividad afecta el procedimiento anterior.

Cuando las restricciones de desigualdad $\text{pags}[S(X)]$ -y se usan en (1.73) en lugar de restricciones de igualdad, el procedimiento de solución sigue siendo casi el mismo. La única diferencia es que el vector multiplicador de Lagrange λ ahora debe ser no negativo. De ello se deduce que el programa dual se convierte en

$$\begin{array}{ll} \text{máximo} & D(\lambda) \\ \lambda & \\ \text{sujeto a:} & \lambda \geq 0, \end{array}$$

con $D(\lambda)$ dado en (1.81).

Una generalización adicional es reemplazar el problema de optimización discreta anterior con un *funcional* problema de optimización. Este tema se discutirá en el Capítulo 8. En particular, la Sección 8.9 trata sobre el método MinxEnt, que involucra un problema funcional de MinxEnt.

PROBLEMAS

Teoría de probabilidad

1.1 Demuestre los siguientes resultados, usando las propiedades de la medida de probabilidad en la Definición 1.2.1 (aquí A y B son eventos):

a) $\text{PAGS}(A^c) = 1 - \text{PAGS}(A)$.

b) $\text{PAGS}(A \cup B) = \text{PAGS}(A) + \text{PAGS}(B) - \text{PAGS}(A \cap B)$.

1.2 Demuestre la regla del producto (1.4) para el caso de tres eventos.

1.3 Sacamos tres bolas consecutivamente de un bol que contiene exactamente cinco bolas blancas y cinco negras, sin volver a colocarlas. ¿Cuál es la probabilidad de que todas las bolas extraídas sean negras?

1.4 Considere el experimento aleatorio en el que lanzamos una moneda sesgada hasta que sale cara. Suponga que la probabilidad de cara en cualquier lanzamiento es pags . Dejar X sea el número de lanzamientos necesarios. Muestra esa X -GRAMO(pags).

1.5 En una sala con mucha gente, le preguntamos a cada persona su cumpleaños (día y mes). Dejar $norte$ sea el número de personas consultadas hasta que obtengamos un cumpleaños "duplicado".

a) Calcular $\text{PAGS}(norte > norte), norte = 0, 1, 2, \dots$

b) Para cual $norte$ tenemos $\text{PAGS}(norte - norte) = 1/2$?

c) Usa una computadora para calcular $MI[norte]$.

1.6 Dejar X y Y sean variables aleatorias normales estándar independientes, y sean tu y V ser variables aleatorias que se derivan de X y Y mediante la transformada lineal

$$\begin{pmatrix} tu \\ V \end{pmatrix} = \begin{pmatrix} \text{pecadoa} & - \text{porquea} \\ \text{porquea} & \text{pecadoa} \end{pmatrix} \begin{pmatrix} X \\ Y \end{pmatrix}.$$

a) Derive el pdf conjunto de tu y V .

b) Muestra esa tu y V son independientes y normalmente distribuidos estándar.

1.7 Dejar $X \sim \text{Exp}(\lambda)$. Muestre que el *propiedad sin memoria* tiene: para todos $s \geq 0$,

$$\text{PAGS}(X > t + s | X > t) = \text{PAGS}(X > s).$$

1.8 Dejar X_1, X_2, X_3 sean variables aleatorias de Bernoulli independientes con probabilidades de éxito $1/2, 1/3$ y $1/4$, respectivamente. Dé su pdf conjunta condicional, dado que $X_1 + X_2 + X_3 = 2$.

1.9 Verifique las expectativas y las variaciones en la Tabla 1.3.

1.10 Dejar X y Y tener densidad articular f dada por

$$f(x, y) = cxy, 0 \leq x \leq 1, 0 \leq y \leq 1.$$

a) Determinar la constante de normalización C .

b) Determinar $\text{PAGS}(X + 2Y - 1)$.

1.11 Dejar $X \sim \text{Exp}(\lambda)$ y $Y \sim \text{Exp}(m)$ Se independiente. Muestra esa

a) $\min(X, Y) \sim \text{Exp}(\lambda + m)$,

b) $\text{PAGS}(X < Y / \min(X, Y)) = \frac{\lambda}{\lambda + m}$.

1.12 Verifique las propiedades de varianza y covarianza en la tabla 1.4.

1.13 Demuestre que el coeficiente de correlación siempre se encuentra entre -1 y 1 . [Sugerencia: utilice el hecho de que la varianza de $hacha + Y$ siempre es no negativo, para cualquier a .]

1.14 Considere los ejemplos 1.1 y 1.2. Definir X como la función que asigna el número $X_1 + \dots + X_{norte}$ a cada resultado $\omega = (X_1, \dots, X_{norte})$. El evento de que haya exactamente k cabezas adentro $norte$ los lanzamientos se pueden escribir como

$$\{\omega \in \Omega : X(\omega) = k\}.$$

Si lo abreviamos a $\{X = k\}$, y abreviar aún más $\text{PAGS}(\{X = k\})$ a $\text{PAGS}(X = k)$, entonces obtenemos exactamente (1.7). Verifique que siempre se pueden ver las variables aleatorias de esta manera, es decir, como funciones de valor real en Ω , y que las probabilidades como $\text{PAGS}(X = X)$ debe interpretarse como $\text{PAGS}(\{\omega \in \Omega : X(\omega) = X\})$.

1.15 Muestra esa

$$\text{Var} \left(\sum_{i=1}^{norte} X_i \right) = \sum_{i=1}^{norte} \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq norte} \text{cov}(X_i, X_j).$$

1.16 Sea Σ la matriz de covarianza de un vector columna aleatorio \mathbf{X} . Escribe $\mathbf{Y} = \mathbf{X} - \mathbf{m}$, dónde \mathbf{m} es el vector de esperanza de \mathbf{X} . Por lo tanto $\Sigma = \text{MI}[\mathbf{AA}]$. Demuestre que Σ es semidefinido positivo. Es decir, para cualquier vector \mathbf{u} , tenemos $\mathbf{u}^T \Sigma \mathbf{u} \geq 0$.

1.17 Suponer $Y \sim \text{Gama}(norte, \lambda)$. Demuestra que para todos $X > 0$

$$\text{PAGS}(Y < X) = 1 - \sum_{k=0}^{norte-1} \frac{e^{-\lambda} \lambda^k}{k!} \frac{\Gamma(norte, \lambda X)}{\Gamma(norte)} . \quad (1.84)$$

1.18 Considere el experimento aleatorio donde dibujamos de manera uniforme e independiente $norte$ números X_1, \dots, X_{norte} , del intervalo $[0, 1]$.

a) Dejar $METRO$ ser el más pequeño de los $norte$ números. Expresar $METRO$ en términos de X_1, \dots, X_{norte} .

b) Determinar el pdf de *METRO*.

1.19 Dejar $Y = \min X$, donde $X \sim \text{NORTE}(0,1)$.

a) Determinar el pdf de Y .

b) Determine el valor esperado de Y .

1.20 Seleccionamos un punto (X, Y) del triángulo $(0,0)-(1,0)-(1,1)$ de tal manera que X tiene una distribución uniforme en $(0,1)$ y la distribución condicional de Y dado $X=x$ es uniforme en $(0, x)$.

a) Determine la función de densidad de probabilidad conjunta de X y Y .

b) Determinar el pdf de Y .

c) Determine la función de densidad de probabilidad condicional de X dado $Y=y$ para todos $y \in (0,1)$.

d) Calcular $MI[X / Y=y]$ para todos $y \in (0,1)$.

mi) Determinar las expectativas de X y Y .

Procesos de veneno

1.21 Dejar $\{norte_t, t=0\}$ sea un proceso de Poisson con tasa $\lambda = 2$. Encuentra

a) $PAGS(norte_2=1, norte_3=4, norte_5=5)$,

b) $PAGS(norte_4=3 / N_2=1, norte_3=2)$,

c) $MI[norte_4 / norte_2=2]$,

d) $PAGS(norte[2,7] = 4, norte[3,8] = 6)$,

mi) $MI[norte[4,6] / norte[1,5] = 3]$.

1.22 Demostrar que para cualquier fijo $k \in \text{NORTE}$, $t > 0$ y $\lambda > 0$,

$$\lim_{norte \rightarrow \infty} \binom{norte}{k} \frac{(\lambda t)^k}{norte^k} \left(1 - \frac{\lambda t}{norte}\right)^{norte-k} = \frac{(\lambda t)^k}{k!} e^{-\lambda t}.$$

[Sugerencia: escriba el coeficiente binomial y use el hecho de que $\lim_{norte \rightarrow \infty} \binom{norte}{k} \left(1 - \frac{\lambda t}{norte}\right)^{norte-k} = e^{-\lambda t}$.]

1.23 Considere la aproximación de Bernoulli en la Sección 1.12. Dejar tu_1, tu_2, \dots denotan los tiempos de éxito para el proceso de Bernoulli X .

a) Verifique que los tiempos de "interéxito" $tu_1, tu_2 - tu_1, \dots$ son independientes y tienen una distribución geométrica con parámetro $p = \lambda h$.

b) Para pequeños h y $norte = t/h$, demostrar que la relación $PAGS(A_1 > t) \approx PAGS(tu_1 > norte)$ lleva en el límite, como $norte \rightarrow \infty$, a

$$PAGS(A_1 > t) = e^{-\lambda t}.$$

1.24 Si $\{norte_t, t=0\}$ es un proceso de Poisson con tasa λ , demuestre que para $0 < tu < ty$, $j=0,1,2, \dots, norte$,

$$PAGS(norte_{tu}=j / N=norte) = \binom{norte}{j} \frac{(\lambda t)^j}{j!} \left(1 - \frac{\lambda t}{norte}\right)^{norte-j},$$

es decir, la distribución condicional de $norte_{tu}$ dado $norte=N=norte$ es binomial con parámetros $norte$ y $\lambda t/norte$.

Procesos de Markov

1.25 Determine el pdf (discreto) de cada X_{norte} , $norte=0,1,2, \dots$, para el paseo aleatorio del ejemplo 1.10. Además, calcule $MI[X_{norte}]$ y la varianza de X_{norte} para cada $norte$.

1.26 Dejar $\{X_{norte}, norte \in norte\}$ ser una cadena de Markov con espacio de estado $\{0,1,2\}$, matriz de transición

$$PAGS = \begin{pmatrix} 0.3 & 0.1 & 0.6 \\ 0.4 & 0.4 & 0.2 \\ 0.1 & 0.7 & 0.2 \end{pmatrix},$$

y distribución inicial $\pi = (0.2, 0.5, 0.3)$. Determinar

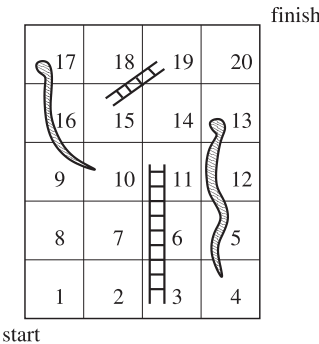
- a) $PAGS(X_1=2)$,
- b) $PAGS(X_2=2)$,
- c) $PAGS(X_3=2 / X_0=0)$,
- d) $PAGS(X_0=1 / X_1=2)$,
- mi) $PAGS(X_1=1, X_3=1)$.

1.27 Dos perros albergan un número total de *metro* Florida fácil Spot inicialmente tiene *b* Floridaeas y Lassie tiene el resto *m-b*. Las pulgas han acordado la siguiente política de inmigración: en todo momento $norte=1,2, \dots$, se selecciona al azar una pulga de la población total y esa pulga saltará de un perro a otro. Describe la población de pulgas en Spot como una cadena de Markov y encuentra su distribución estacionaria.

1.28 Clasifique los estados de la cadena de Markov con la siguiente matriz de transición:

$$PAGS = \begin{pmatrix} 0.0 & 0.3 & 0.6 & 0.0 & 0.1 \\ 0.0 & 0.3 & 0.0 & 0.7 & 0.0 \\ 0.3 & 0.1 & 0.6 & 0.0 & 0.0 \\ 0.0 & 0.1 & 0.0 & 0.9 & 0.0 \\ 0.1 & 0.1 & 0.2 & 0.0 & 0.6 \end{pmatrix}.$$

1.29 Considere el siguiente juego de serpientes y escaleras. Dejar $norte$ sea el número de lanzamientos requeridos para llegar al final usando un dado justo. Calcular la esperanza de $norte$ usando una computadora.



1.30 La Sra. Ella Brum camina de un lado a otro entre su casa y su oficina todos los días. Tiene tres paraguas, que se distribuyen en dos paragueros (uno en casa y otro en el trabajo). Cuando no llueve, la Sra. Brum camina sin paraguas. Cuando llueve, toma un paraguas del puesto en el lugar de

su partida, siempre que haya uno disponible. Suponga que la probabilidad de que esté lloviendo en el momento de cualquier salida es p . Dejar X_{norte} indique el número de sombrillas disponibles en el lugar donde Ella llega después del número de caminata $norte$; $norte=1, 2, \dots$, incluido el que posiblemente traiga consigo. Calcule la probabilidad límite de que llueva y no haya paraguas disponible.

1.31 Se suelta un ratón en el laberinto de la figura 1.10. De cada compartimento, el ratón elige uno de los compartimentos adyacentes con igual probabilidad, independientemente del pasado. El ratón pasa una cantidad de tiempo distribuida exponencialmente en cada compartimento. El tiempo medio de permanencia en cada uno de los compartimentos 1, 3 y 4 es de dos segundos; el tiempo medio de permanencia en los compartimentos 2, 5 y 6 es de cuatro segundos. Dejar $\{X_t, t \geq 0\}$ sea el proceso de salto de Markov que describe la posición del ratón por veces $t \geq 0$. Suponga que el mouse comienza en el compartimento 1 en el momento $t=0$.

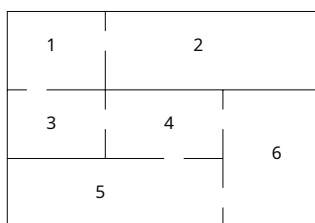


Figura 1.10: Un laberinto.

¿Cuáles son las probabilidades de que el ratón se encuentre en cada uno de los compartimentos 1, 2, \dots , 6 en algún momento t muy lejos en el futuro?

1.32 En un $M/M/\infty$ -sistema de colas, los clientes llegan de acuerdo con un proceso de Poisson con tasa a . Cada cliente que ingresa es inmediatamente atendido por uno de una infinidad de servidores; por lo tanto, no hay cola. Los tiempos de servicio están distribuidos exponencialmente, con media $1/b$. Todos los tiempos de servicio y entre llegadas son independientes. Dejar X_t sea el número de clientes en el sistema en ese momento t . Demuestre que la distribución límite de X_t , como $t \rightarrow \infty$, es Poisson con parámetro un/b .

Mejoramiento

1.33 Dejar \mathbf{a} y \mathbf{X} ser $norte$ -vectores columna dimensionales. Muestra que $\mathbf{X}^T \mathbf{a} = \mathbf{a}^T \mathbf{X}$.

1.34 Dejar \mathbf{A} ser una simétrica $norte \times norte$ matriz y \mathbf{X} un $norte$ -vector columna dimensional. Muestra que $\mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{A} \mathbf{X}^T \mathbf{X}$. ¿Cuál es el gradiente si \mathbf{A} no es simétrico?

1.35 Demuestre que la distribución óptima \mathbf{p}^* en el ejemplo 1.16 está dada por la distribución uniforme.

1.36 Derive el programa (1.71).

1.37 Considere el programa MinxEnt

$$\begin{aligned} \min_{\mathbf{pags}} \quad & \sum_{i=1}^{norte} \frac{pags_i}{q_i} \\ \text{sujeto a:} \quad & \mathbf{pags} \geq \mathbf{0}, \quad A\mathbf{pags} = \mathbf{b}, \quad \sum_{i=1}^{norte} pags_i = 1, \end{aligned}$$

dónde \mathbf{pags} y \mathbf{q} son vectores de distribución de probabilidad y A es un $norte \times norte$ matriz.

a) Demuestre que el Lagrangiano para este problema es de la forma

$$L(\mathbf{pags}, \lambda, \beta, \mathbf{m}) = \mathbf{pags}^T \mathbf{q} - \lambda(A\mathbf{pags} - \mathbf{b}) - \mathbf{m}^T (\mathbf{pags} - \mathbf{1}).$$

b) Muestra que $pags_i = q_i \exp(-\beta - 1 + \mathbf{m}^T \sum_{j=1}^{norte} \lambda_j a_{ji})$, por $i=1, \dots, norte$.

c) Explique por qué, como resultado de las condiciones KKT, el óptimo \mathbf{m} debe ser igual al vector cero.

d) Demuestre que la solución para este programa MinxEnt es exactamente la misma que para el programa donde se omiten las restricciones de no negatividad.

Otras lecturas

Una introducción fácil a la teoría de la probabilidad con muchos ejemplos es [13], y un libro de texto más detallado es [8]. Una referencia clásica es [6]. En [3] se da un tratamiento preciso y accesible de varios procesos estocásticos. Para la optimización convexa nos referimos a [2] y [7].

REFERENCIAS

1. ZI Botev, DP Kroese y T. Taimre. Métodos generalizados de entropía cruzada para la simulación y optimización de eventos raros. *Simulación: transacciones de la Society for Modeling and Simulation International*, 83(11):785–806, 2007.
2. S. Boyd y L. Vandenberghe. *Optimización convexa*. Prensa de la Universidad de Cambridge, Cambridge, Reino Unido, 2004.
3. E. Çinlar. *Introducción a los Procesos Estocásticos*. Prentice Hall, Englewood Cliffs, Nueva Jersey, 1975.
4. TM Cover y JA Thomas. *Elementos de la teoría de la información*. John Wiley & Sons, Nueva York, 1991.
5. CW Curtis. *Álgebra lineal: un enfoque introductorio*. Springer-Verlag, Nueva York, 1984.
6. W. Feller. *Una introducción a la teoría de la probabilidad y sus aplicaciones*, volumen 1. John Wiley & Sons, Nueva York, 2ª edición, 1970.
7. R. Fletcher. *Métodos prácticos de optimización*. John Wiley & Sons, Nueva York, 1987.
8. GR Grimmett y DR Stirzaker. *Procesos de probabilidad y aleatorios*. Oxford University Press, Oxford, 3ª edición, 2001.
9. JN Kapur y HK Kesavan. *Principios de optimización de entropía con aplicaciones*. Prensa académica, Nueva York, 1992.