

Dirk P. Kroese · Joshua C.C. Chan

Statistical Modeling and Computation



Springer

Statistical Modeling and Computation

Dirk P. Kroese • Joshua C.C. Chan

Statistical Modeling and Computation



Springer

Dirk P. Kroese
The University of Queensland
School of Mathematics and Physics
Brisbane, Australia

Joshua C.C. Chan
Department of Economics
Australian National University
Canberra, Australia

ISBN 978-1-4614-8774-6 ISBN 978-1-4614-8775-3 (eBook)

DOI 10.1007/978-1-4614-8775-3

Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2013948920

© The Author(s) 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

*In memory of Reuven Rubinstein, my Friend
and Mentor*


Dirk Kroese

To Raquel

Joshua Chan

Preface

Statistics provides one of the few principled means to extract information from random data and has perhaps more interdisciplinary connections than any other field of science. However, for a beginning student of statistics, the abundance of mathematical concepts, statistical philosophies, and numerical techniques can seem overwhelming. The purpose of this book is to provide a comprehensive and accessible introduction to modern statistics, illuminating its many facets, from both classical (frequentist) and Bayesian points of view. The book offers an integrated treatment of mathematical statistics and modern statistical computation.

The book is aimed at beginning students of statistics and practitioners who would like to fully understand the theory and key numerical techniques of statistics. It is based on a progression of undergraduate statistics courses at The University of Queensland and the Australian National University. Parts of the book have also been successfully tested at the University of New South Wales. Emphasis is laid on the mathematical and computational aspects of statistics. No prior knowledge of statistics is required, but we assume that the reader has a basic knowledge of mathematics, which forms an essential basis for the development of the statistical theory. Starting from scratch, the book gradually builds up to an advanced undergraduate level, providing a solid basis for possible postgraduate research. Throughout the text we illustrate the theory by providing working code in MATLAB, rather than relying on black-box statistical packages. We make frequent use of the symbol  in the margin to facilitate cross-referencing between related pages. The book is accompanied by the web site www.statmodcomp.org from which the MATLAB code and data files can be downloaded. In addition, we provide an R equivalent for each MATLAB program.

The book is structured into three parts. In Part I we introduce the fundamentals of probability theory. We discuss models for random experiments, conditional probability and independence, random variables, and probability distributions. Moreover, we explain how to carry out random experiments on a computer.

In Part II we introduce the general framework for statistical modeling and inference, from both classical and Bayesian perspectives. We discuss a variety of common models for data, such as independent random samples, linear regression,

and ANOVA models. Once a model for the data is determined one can carry out a mathematical analysis of the model on the basis of the available data. We discuss a wide range of concepts and techniques for statistical inference, including likelihood-based estimation and hypothesis testing, sufficiency, confidence intervals, and kernel density estimation. We encompass both classical and Bayesian approaches and also highlight popular Monte Carlo sampling techniques.

In Part III we address the statistical analysis and computation of a variety of advanced models, such as generalized linear models, autoregressive and moving average models, Gaussian models, and state space models. Particular attention is paid to fast numerical techniques for classical and Bayesian inference on these models. Throughout the book our leading principle is that the mathematical formulation of a statistical model goes hand in hand with the specification of its simulation counterpart.

The book contains a large number of illustrative examples and problem sets (with solutions). To keep the book fully self-contained, we include the more technical proofs and mathematical theory in Appendix B. Appendix A features a concise introduction to MATLAB.

Brisbane, Australia
Canberra, Australia

Dirk Kroese
Joshua Chan

Acknowledgements

This book has benefited from the input of many people. We thank Zdravko Botev, Tim Brereton, Hyun Choi, Eric Eisenstat, Eunice Foo, Catherine Forbes, Patricia Galvan, Ivan Jeliazkov, Ross McVinish, Gary Koop, Rongrong Qu, Ad Ridder, Leonardo Rojas–Nandayapa, John Stachurski, Rodney Strachan, Mingzhu Sun, Thomas Taimre, Justin Tobias, Elisse Yulian, and Bo Zhang for their valuable comments and suggestions on previous drafts of the book.

Contents

Part I Fundamentals of Probability

1	Probability Models	3
1.1	Random Experiments	3
1.2	Sample Space	5
1.3	Events	6
1.4	Probability	9
1.5	Conditional Probability and Independence	12
1.5.1	Product Rule	14
1.5.2	Law of Total Probability and Bayes' Rule	16
1.5.3	Independence	17
1.6	Problems	19
2	Random Variables and Probability Distributions	23
2.1	Random Variables	23
2.2	Probability Distribution	25
2.2.1	Discrete Distributions	27
2.2.2	Continuous Distributions	28
2.3	Expectation	29
2.4	Transforms	33
2.5	Common Discrete Distributions	36
2.5.1	Bernoulli Distribution	36
2.5.2	Binomial Distribution	37
2.5.3	Geometric Distribution	38
2.5.4	Poisson Distribution	40
2.6	Common Continuous Distributions	42
2.6.1	Uniform Distribution	42
2.6.2	Exponential Distribution	43
2.6.3	Normal (Gaussian) Distribution	45
2.6.4	Gamma and χ^2 Distribution	48

2.6.5	<i>F</i> Distribution	49
2.6.6	Student's <i>t</i> Distribution	50
2.7	Generating Random Variables	51
2.7.1	Generating Uniform Random Variables	52
2.7.2	Inverse-Transform Method	53
2.7.3	Acceptance–Rejection Method	55
2.8	Problems	57
3	Joint Distributions	63
3.1	Discrete Joint Distributions	64
3.1.1	Multinomial Distribution	68
3.2	Continuous Joint Distributions	69
3.3	Mixed Joint Distributions	73
3.4	Expectations for Joint Distributions	74
3.5	Functions of Random Variables	78
3.5.1	Linear Transformations	79
3.5.2	General Transformations	81
3.6	Multivariate Normal Distribution	82
3.7	Limit Theorems	89
3.8	Problems	93
 Part II Statistical Modeling and Classical and Bayesian Inference		
4	Common Statistical Models	101
4.1	Independent Sampling from a Fixed Distribution	101
4.2	Multiple Independent Samples	103
4.3	Regression Models	104
4.3.1	Simple Linear Regression	105
4.3.2	Multiple Linear Regression	106
4.3.3	Regression in General	108
4.4	Analysis of Variance Models	111
4.4.1	Single-Factor ANOVA	111
4.4.2	Two-Factor ANOVA	113
4.5	Normal Linear Model	114
4.6	Problems	118
5	Statistical Inference	121
5.1	Estimation	122
5.1.1	Method of Moments	123
5.1.2	Least-Squares Estimation	125
5.2	Confidence Intervals	128
5.2.1	Iid Data: Approximate Confidence Interval for μ	130
5.2.2	Normal Data: Confidence Intervals for μ and σ^2	131

5.2.3	Two Normal Samples: Confidence Intervals for $\mu_X - \mu_Y$ and σ_X^2/σ_Y^2	133
5.2.4	Binomial Data: Approximate Confidence Intervals for Proportions	135
5.2.5	Confidence Intervals for the Normal Linear Model	137
5.3	Hypothesis Testing	140
5.3.1	ANOVA for the Normal Linear Model	142
5.4	Cross-Validation	146
5.5	Sufficiency and Exponential Families	150
5.6	Problems	154
6	Likelihood	161
6.1	Log-Likelihood and Score Functions	165
6.2	Fisher Information and Cramér–Rao Inequality	167
6.3	Likelihood Methods for Estimation	172
6.3.1	Score Intervals	174
6.3.2	Properties of the ML Estimator	175
6.4	Likelihood Methods in Statistical Tests	178
6.5	Newton–Raphson Method	180
6.6	Expectation–Maximization (EM) Algorithm	182
6.7	Problems	188
7	Monte Carlo Sampling	195
7.1	Empirical Cdf	196
7.2	Density Estimation	201
7.3	Resampling and the Bootstrap Method	203
7.4	Markov Chain Monte Carlo	209
7.5	Metropolis–Hastings Algorithm	214
7.6	Gibbs Sampler	218
7.7	Problems	220
8	Bayesian Inference	227
8.1	Hierarchical Bayesian Models	229
8.2	Common Bayesian Models	233
8.2.1	Normal Model with Unknown μ and σ^2	233
8.2.2	Bayesian Normal Linear Model	237
8.2.3	Bayesian Multinomial Model	240
8.3	Bayesian Networks	244
8.4	Asymptotic Normality of the Posterior Distribution	248
8.5	Priors and Conjugacy	249
8.6	Bayesian Model Comparison	251
8.7	Problems	256
 Part III Advanced Models and Inference		
9	Generalized Linear Models	265
9.1	Generalized Linear Models	265

9.2	Logit and Probit Models	267
9.2.1	Logit Model	267
9.2.2	Probit Model	273
9.2.3	Latent Variable Representation	278
9.3	Poisson Regression	282
9.4	Problems	284
10	Dependent Data Models	287
10.1	Autoregressive and Moving Average Models	287
10.1.1	Autoregressive Models	287
10.1.2	Moving Average Models	297
10.1.3	Autoregressive-Moving Average Models	303
10.2	Gaussian Models	305
10.2.1	Gaussian Graphical Model	306
10.2.2	Random Effects	308
10.2.3	Gaussian Linear Mixed Models	315
10.3	Problems	320
11	State Space Models	323
11.1	Unobserved Components Model	325
11.1.1	Classical Estimation	327
11.1.2	Bayesian Estimation	332
11.2	Time-Varying Parameter Model	333
11.2.1	Bayesian Estimation	334
11.3	Stochastic Volatility Model	339
11.3.1	Auxiliary Mixture Sampling Approach	340
11.4	Problems	346
A	Matlab Primer	349
A.1	Matrices and Matrix Operations	349
A.2	Some Useful Built-In Functions	352
A.3	Flow Control	354
A.4	Function Handles and Function Files	355
A.5	Graphics	356
A.6	Optimization Routines	360
A.7	Handling Sparse Matrices	362
A.8	Gamma and Dirichlet Generator	364
A.9	Cdfs and Inverse Cdfs	365
A.10	Further Reading and References	366
B	Mathematical Supplement	367
B.1	Multivariate Differentiation	367
B.2	Proof of Theorem 2.6 and Corollary 2.2	369
B.3	Proof of Theorem 2.7	370
B.4	Proof of Theorem 3.10	371
B.5	Proof of Theorem 5.2	371

Contents	xv
References	373
Solutions	375
Index	395

Abbreviations and Acronyms

ANOVA	Analysis of variance
AR	Autoregressive
ARMA	Autoregressive-moving average
cdf	Cumulative distribution function
EM	Expectation–maximization
iid	Independent and identically distributed
pdf	Probability density function (discrete or continuous)
PGF	Probability generating function
KDE	Kernel density estimate/estimator
MA	Moving average
MCMC	Markov chain Monte Carlo
MGF	Moment generating function
ML(E)	Maximum likelihood (estimate/estimator)
PRESS	Predicted residual sum of squares

Mathematical Notation

Throughout this book we use notation in which different fonts and letter cases signify different types of mathematical objects. For example, vectors $\mathbf{a}, \mathbf{b}, \mathbf{x}, \dots$ are written in lowercase boldface font and matrices A, B, X in uppercase normal font. Sans serif fonts indicate probability distributions, such as N , Exp , and Bin . Probability and expectation symbols are written in blackboard bold font: \mathbb{P} and \mathbb{E} . MATLAB code and functions will always be written in typewriter font.

Traditionally, classical and Bayesian statistics use a *different* notation system for random variables and their probability density functions. In classical statistics and probability theory random variables usually are denoted by uppercase letters X, Y, Z, \dots and their outcomes by lowercase letters x, y, z, \dots . Bayesian statisticians typically use lowercase letters for both. More importantly, in the Bayesian notation system, it is common to use the *same* letter f (or p) for different probability densities, as in $f(x, y) = f(x)f(y)$. Classical statisticians and probabilists would prefer a different symbol for each function, as in $f(x, y) = f_X(x)f_Y(y)$. We will predominantly use the classical notation, especially in the first part of the book. However, when dealing with Bayesian models and inference, such as in Chaps. 8 and 11, it will be convenient to switch to the Bayesian notation system. Here is a list of frequently used symbols:

\approx	Is approximately
\propto	Is proportional to
∞	Infinity
\otimes	Kronecker product
$\stackrel{\text{def}}{=}$	Is defined as
\sim	Is distributed as
$\overset{\text{iid}}{\sim}, \sim_{\text{iid}}$	Are independent and identically distributed as
$\overset{\text{approx.}}{\sim}$	Is approximately distributed as
\mapsto	Maps to
$A \cup B$	Union of sets A and B
$A \cap B$	Intersection of sets A and B
A^c	Complement of set A
$A \subset B$	A is a subset of B
\emptyset	Empty set
$\ \mathbf{x}\ $	Euclidean norm of vector \mathbf{x}
∇f	Gradient of f
$\nabla^2 f$	Hessian of f
A^\top, \mathbf{x}^\top	Transpose of matrix A or vector \mathbf{x}
$\text{diag}(\mathbf{a})$	Diagonal matrix with diagonal entries defined by \mathbf{a}
$\text{tr}(A)$	Trace of matrix A
$\det(A)$	Determinant of matrix A

$ A $	Absolute value of the determinant of matrix A . Also, number of elements in set A or absolute value of real number A
argmax	$\operatorname{argmax} f(x)$ is a value x^* for which $f(x^*) \geq f(x)$ for all x
d	Differential symbol
\mathbb{E}	Expectation
e	Euler's constant $\lim_{n \rightarrow \infty} (1 + 1/n)^n = 2.71828 \dots$
$I_A, I\{A\}$	Indicator function: equal to 1 if the condition/event A holds and 0 otherwise.
\ln	(Natural) logarithm
\mathbb{N}	Set of natural numbers $\{0, 1, \dots\}$
φ	Pdf of the standard normal distribution
Φ	Cdf of the standard normal distribution
\mathbb{P}	Probability measure
\mathcal{O}	Big-O order symbol: $f(x) = \mathcal{O}(g(x))$ if $ f(x) \leq \alpha g(x)$ for some constant α as $x \rightarrow a$
o	Little-o order symbol: $f(x) = o(g(x))$ if $f(x)/g(x) \rightarrow 0$ as $x \rightarrow a$
\mathbb{R}	The real line = one-dimensional Euclidean space
\mathbb{R}_+	Positive real line: $[0, \infty)$
\mathbb{R}^n	n -Dimensional Euclidean space
$\hat{\theta}$	Estimate/estimator
\mathbf{x}, \mathbf{y}	Vectors
\mathbf{X}, \mathbf{Y}	Random vectors
\mathbb{Z}	Set of integers $\{\dots, -1, 0, 1, \dots\}$

Probability Distributions

Ber	Bernoulli distribution
Beta	Beta distribution
Bin	Binomial distribution
Cauchy	Cauchy distribution
χ^2	Chi-squared distribution
Dirichlet	Dirichlet distribution
DU	Discrete uniform distribution
Exp	Exponential distribution
F	F distribution
Gamma	Gamma distribution
Geom	Geometric distribution
InvGamma	Inverse-gamma distribution
Mnom	Multinomial distribution
N	Normal or Gaussian distribution
Poi	Poisson distribution
t	Student's t distribution
TN	Truncated normal distribution
U	Uniform distribution
Weib	Weibull distribution

Part I

Fundamentals of Probability

In Part I of the book we consider the *probability* side of statistics. In particular, we will consider how random experiments can be modeled mathematically and how such modeling enables us to compute various properties of interest for those experiments.

Chapter 1

Probability Models

1.1 Random Experiments

The basic notion in probability is that of a **random experiment**: an experiment whose outcome cannot be determined in advance, but which is nevertheless subject to analysis. Examples of random experiments are:

1. Tossing a die and observing its face value.
2. Measuring the amount of monthly rainfall in a certain location.
3. Counting the number of calls arriving at a telephone exchange during a fixed time period.
4. Selecting at random fifty people and observing the number of left-handers.
5. Choosing at random ten people and measuring their heights.

The goal of *probability* is to understand the behavior of random experiments by analyzing the corresponding *mathematical models*. Given a mathematical model for a random experiment one can calculate quantities of interest such as probabilities and expectations. Moreover, such mathematical models can typically be implemented on a computer, so that it becomes possible to *simulate* the experiment. Conversely, any computer implementation of a random experiment implicitly defines a mathematical model. Mathematical models for random experiments are also the basis of *statistics*, where the objective is to infer which of several competing models best fits the observed data. This often involves the estimation of model parameters from the data.

Example 1.1 (Coin Tossing). One of the most fundamental random experiments is the one where a coin is tossed a number of times. Indeed, much of probability theory can be based on this simple experiment. To better understand how this coin toss experiment behaves, we can carry it out on a computer, using programs such as MATLAB. The following simple MATLAB program simulates a sequence of 100 tosses with a fair coin (i.e., Heads and Tails are equally likely) and plots the results in a bar chart.

```
x = (rand(1,100) < 0.5)    % generate the coin tosses
bar(x)                     % plot the results in a bar chart
```

The function `rand` draws uniform random numbers from the interval $[0, 1]$ —in this case a 1×100 vector of such numbers. By testing whether the uniform numbers are less than 0.5, we obtain a vector `x` of 1s and 0s, indicating Heads and Tails, say. Typical outcomes for three such experiments are given in Fig. 1.1.

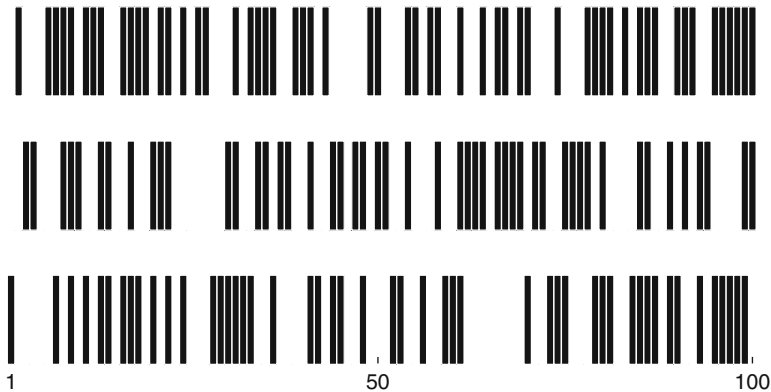


Fig. 1.1 Three experiments where a fair coin is tossed 100 times. The *dark bars* indicate when “Heads” (= 1) appears

We can also plot the average number of Heads against the number of tosses. In the same MATLAB program, this is accomplished by adding two lines of code:

```
y = cumsum(x) ./ [1:100] % calculate the cumulative sum and
                        % divide this elementwise by the vector [1:100]
plot(y)                % plot the result in a line graph
```

The result of three such experiments is depicted in Fig. 1.2. Notice that the average number of Heads seems to converge to 0.5, but there is a lot of random fluctuation.

Similar results can be obtained for the case where the coin is *biased*, with a probability of Heads of p , say. Here are some typical *probability* questions.

- What is the probability of x Heads in 100 tosses?
- What is the expected number of Heads?
- How long does one have to wait until the first Head is tossed?
- How fast does the average number of Heads converge to p ?

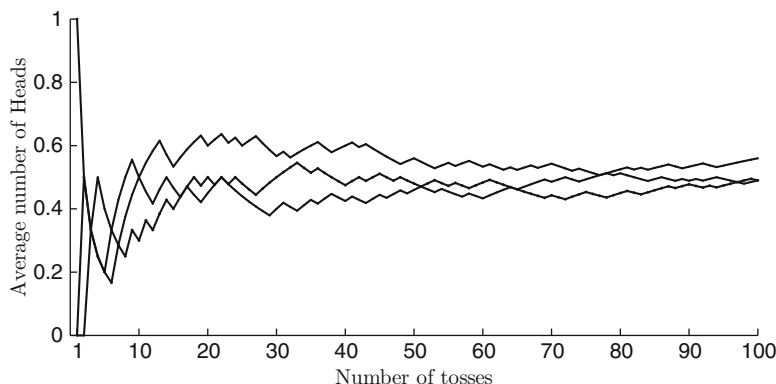


Fig. 1.2 The average number of Heads in n tosses, where $n = 1, \dots, 100$

A statistical analysis would start from observed data of the experiment—for example, all the outcomes of 100 tosses are known. Suppose the probability of Heads p is not known. Typical *statistics* questions are:

- Is the coin fair?
- How can p be best estimated from the data?
- How accurate/reliable would such an estimate be?

The mathematical models that are used to describe random experiments consist of three building blocks: a *sample space*, a set of *events*, and a *probability*. We will now describe each of these objects.

1.2 Sample Space

Although we cannot predict the outcome of a random experiment with certainty, we usually can specify a set of possible outcomes. This gives the first ingredient in our model for a random experiment.

Definition 1.1. (Sample Space). The **sample space** Ω of a random experiment is the set of all possible outcomes of the experiment.

Examples of random experiments with their sample spaces are:

1. Cast two dice consecutively and observe their face values:

$$\Omega = \{(1, 1), (1, 2), \dots, (1, 6), (2, 1), \dots, (6, 6)\}.$$

2. Measure the lifetime of a machine in days:

$$\Omega = \mathbb{R}_+ = \{ \text{positive real numbers} \}.$$

3. Count the number of arriving calls at an exchange during a specified time interval:

$$\Omega = \{0, 1, \dots\}.$$

4. Measure the heights of 10 people:

$$\Omega = \{(x_1, \dots, x_{10}) : x_i \geq 0, i = 1, \dots, 10\} = \mathbb{R}_+^{10}.$$

Here (x_1, \dots, x_{10}) represents the outcome that the height of the first selected person is x_1 , the height of the second person is x_2 , and so on.

Notice that for modeling purposes it is often easier to take the sample space larger than is strictly necessary. For example, the actual lifetime of a machine would in reality not span the entire positive real axis, and the heights of the ten selected people would not exceed 9 ft.

1.3 Events

Often we are not interested in a single outcome but in whether or not one in a *group* of outcomes occurs.

Definition 1.2. (Event). An **event** is a subset of the sample space Ω to which a probability can be assigned.

Events will be denoted by capital letters A, B, C, \dots . We say that event A **occurs** if the outcome of the experiment is one of the elements in A .

Examples of events are:

1. The event that the sum of two dice is 10 or more:

$$A = \{(4, 6), (5, 5), (5, 6), (6, 4), (6, 5), (6, 6)\}.$$

2. The event that a machine is functioning for less than 1000 days:

$$A = [0, 1000).$$

3. The event that out of a group of 50 people 5 are left-handed:

$$A = \{5\}.$$

Example 1.2 (Coin Tossing). Suppose that a coin is tossed 3 times and that we record either Heads or Tails at every toss. The sample space can then be written as

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

where, for instance, HTH means that the first toss is Heads, the second Tails, and the third Heads. An alternative (but equivalent) sample space is the set $\{0, 1\}^3$ of binary vectors of length 3; for example, HTH corresponds to $(1, 0, 1)$ and THH to $(0, 1, 1)$.

The event A that the third toss is Heads is

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

Since events are sets, we can apply the usual set operations to them, as illustrated in the *Venn diagrams* in Fig. 1.3.

1. The set $A \cap B$ (A **intersection** B) is the event that A and B both occur.
2. The set $A \cup B$ (A **union** B) is the event that A or B or both occur.
3. The event A^c (A **complement**) is the event that A does *not* occur.
4. If $B \subset A$ (B is a **subset** of A), then event B is said to *imply* event A .

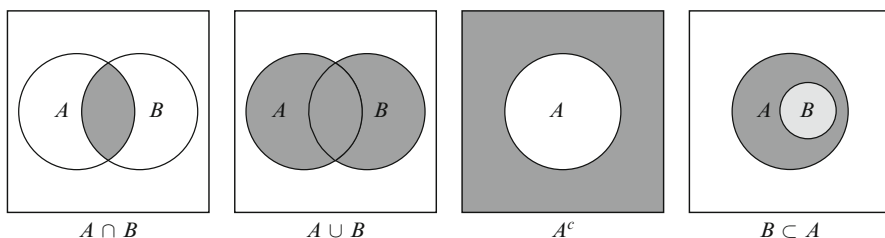


Fig. 1.3 Venn diagrams of set operations. Each *square* represents the sample space Ω

Two events A and B which have no outcomes in common, that is, $A \cap B = \emptyset$ (empty set), are called **disjoint** events.

Example 1.3 (Casting Two Dice). Suppose we cast two dice consecutively. The sample space is $\Omega = \{(1, 1), (1, 2), \dots, (1, 6), (2, 1), \dots, (6, 6)\}$. Let $A = \{(6, 1), \dots, (6, 6)\}$ be the event that the first die is 6, and let $B = \{(1, 6), \dots, (6, 6)\}$ be the event that the second die is 6. Then $A \cap B = \{(6, 1), \dots, (6, 6)\} \cap \{(1, 6), \dots, (6, 6)\} = \{(6, 6)\}$ is the event that both dice are 6.

Example 1.4 (System Reliability). In Fig. 1.4 three systems are depicted, each consisting of three unreliable components. The *series* system works if all components work; the *parallel* system works if at least one of the components works; and the *2-out-of-3* system works if at least 2 out of 3 components work.

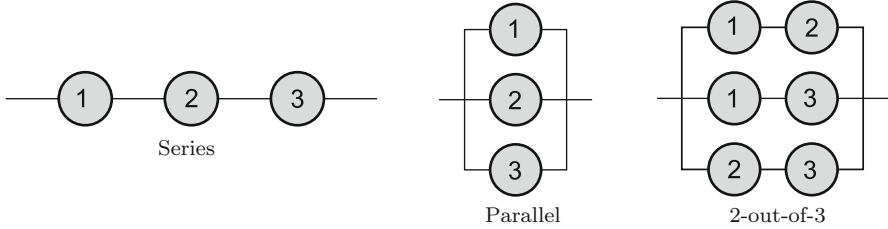


Fig. 1.4 Three unreliable systems

Let A_i be the event that the i th component is functioning, $i = 1, 2, 3$; and let D_a, D_b, D_c be the events that, respectively, the series, parallel, and 2-out-of-3 system are functioning. Then, $D_a = A_1 \cap A_2 \cap A_3$ and $D_b = A_1 \cup A_2 \cup A_3$. Also,

$$\begin{aligned} D_c &= (A_1 \cap A_2 \cap A_3) \cup (A_1^c \cap A_2 \cap A_3) \cup (A_1 \cap A_2^c \cap A_3) \cup (A_1 \cap A_2 \cap A_3^c) \\ &= (A_1 \cap A_2) \cup (A_1 \cap A_3) \cup (A_2 \cap A_3) . \end{aligned}$$

Two useful results in the theory of sets are the following, due to De Morgan:

Theorem 1.1. (De Morgan's Laws). If $\{A_i\}$ is a collection of sets, then

$$\left(\bigcup_i A_i \right)^c = \bigcap_i A_i^c \quad (1.1)$$

and

$$\left(\bigcap_i A_i \right)^c = \bigcup_i A_i^c . \quad (1.2)$$

Proof. If we interpret A_i as the event that component i works in Example 1.4, then the left-hand side of (1.1) is the event that the parallel system is not working. The right-hand side of (1.1) is the event that all components are not working. Clearly these two events are identical. The proof for (1.2) follows from a similar reasoning;

1.4 Probability

The third ingredient in the model for a random experiment is the specification of the probability of the events. It tells us how *likely* it is that a particular event will occur.

Definition 1.3. (Probability). A **probability** \mathbb{P} is a function which assigns a number between 0 and 1 to each event and which satisfies the following rules:

1. $0 \leq \mathbb{P}(A) \leq 1$.
2. $\mathbb{P}(\Omega) = 1$.
3. For any sequence A_1, A_2, \dots of *disjoint* events we have

$$\textbf{Sum Rule:} \quad \mathbb{P}\left(\bigcup_i A_i\right) = \sum_i \mathbb{P}(A_i) . \quad (1.3)$$

The crucial property (1.3) is called the **sum rule** of probability. It simply states that if an event can happen in several distinct ways (expressed as a union of events, none of which are overlapping), then the probability that at least one of these events happens (i.e., the probability of the union) is simply the sum of the probabilities of the individual events. Figure 1.5 illustrates that the probability \mathbb{P} has the properties of a *measure*. However, instead of measuring lengths, areas, or volumes, $\mathbb{P}(A)$ measures the likelihood or probability of an event A as a number between 0 and 1.

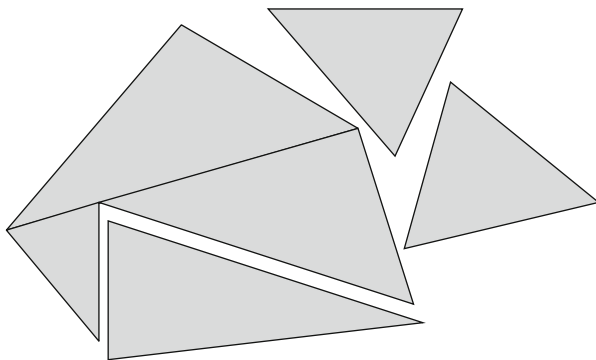


Fig. 1.5 A probability rule \mathbb{P} has exactly the same properties as an area measure. For example, the total area of the union of the nonoverlapping triangles is equal to the sum of the areas of the individual triangles

The following theorem lists some important properties of a probability measure. These properties are direct consequences of the three rules defining a probability measure.

Theorem 1.2. (Properties of a Probability). Let A and B be events and \mathbb{P} a probability. Then,

1. $\mathbb{P}(\emptyset) = 0$,
2. if $A \subset B$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$,
3. $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$,
4. $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$.

Proof.

1. Since $\Omega = \Omega \cup \emptyset$ and $\Omega \cap \emptyset = \emptyset$, it follows from the sum rule that $\mathbb{P}(\Omega) = \mathbb{P}(\Omega) + \mathbb{P}(\emptyset)$. Therefore, by Rule 2 of Definition 1.3, we have $1 = 1 + \mathbb{P}(\emptyset)$, from which it follows that $\mathbb{P}(\emptyset) = 0$.
2. If $A \subset B$, then $B = A \cup (B \cap A^c)$, where A and $B \cap A^c$ are disjoint. Hence, by the sum rule, $\mathbb{P}(B) = \mathbb{P}(A) + \mathbb{P}(B \cap A^c)$, which (by Rule 1) is greater than or equal to $\mathbb{P}(A)$.
3. $\Omega = A \cup A^c$, where A and A^c are disjoint. Hence, by the sum rule and Rule 2: $1 = \mathbb{P}(\Omega) = \mathbb{P}(A) + \mathbb{P}(A^c)$, and thus $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$.
4. Write $A \cup B$ as the disjoint union of A and $B \cap A^c$. Then, $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B \cap A^c)$. Also, $B = (A \cap B) \cup (B \cap A^c)$, so that $\mathbb{P}(B) = \mathbb{P}(A \cap B) + \mathbb{P}(B \cap A^c)$. Combining these two equations gives $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$. \square

We have now completed our general model for a random experiment. Of course for any *specific* model we must carefully specify the sample space Ω and probability \mathbb{P} that best describe the random experiment.

Example 1.5 (Casting a Die). Consider the experiment where a fair die is cast. How should we specify Ω and \mathbb{P} ? Obviously, $\Omega = \{1, 2, \dots, 6\}$; and common sense dictates that we should define \mathbb{P} by

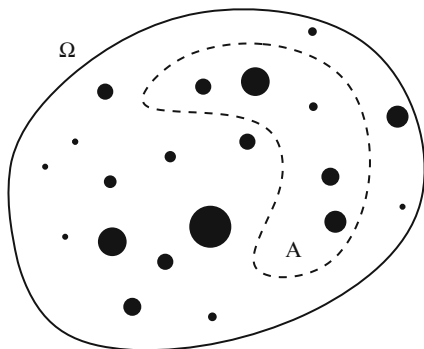
$$\mathbb{P}(A) = \frac{|A|}{6}, \quad A \subset \Omega ,$$

where $|A|$ denotes the number of elements in set A . For example, the probability of getting an even number is $\mathbb{P}(\{2, 4, 6\}) = 3/6 = 1/2$.

In many applications the sample space is *countable*: $\Omega = \{a_1, a_2, \dots, a_n\}$ or $\Omega = \{a_1, a_2, \dots\}$. Such a sample space is said to be **discrete**. The easiest way to specify a probability \mathbb{P} on a discrete sample space is to first assign a probability p_i to each **elementary event** $\{a_i\}$ and then to define

$$\mathbb{P}(A) = \sum_{i: a_i \in A} p_i \quad \text{for all } A \subset \Omega .$$

Fig. 1.6 A discrete sample space



This idea is graphically represented in Fig. 1.6. Each element a_i in the sample space is assigned a probability weight p_i represented by a black dot. To find the probability of an event A we have to sum up the weights of all the elements in the set A .

Again, it is up to the modeler to properly specify these probabilities. Fortunately, in many applications, all elementary events are *equally likely*, and thus the probability of each elementary event is equal to 1 divided by the total number of elements in Ω . In such case the probability of an event $A \subset \Omega$ is simply

$$\mathbb{P}(A) = \frac{|A|}{|\Omega|} = \frac{\text{Number of elements in } A}{\text{Number of elements in } \Omega},$$

provided that the total number of elements in Ω is finite. The calculation of such probabilities thus reduces to *counting*; see Problem 1.6.

19

When the sample space is not countable, for example, $\Omega = \mathbb{R}_+$, it is said to be **continuous**.

Example 1.6 (Drawing a Random Point in the Unit Interval). We draw at random a point in the interval $[0, 1]$ such that each point is equally likely to be drawn. How do we specify the model for this experiment?

The sample space is obviously $\Omega = [0, 1]$, which is a continuous sample space. We cannot define \mathbb{P} via the elementary events $\{x\}$, $x \in [0, 1]$ because each of these events has probability 0. However, we can define \mathbb{P} as follows. For each $0 \leq a \leq b \leq 1$, let

$$\mathbb{P}([a, b]) = b - a.$$

This completely defines \mathbb{P} . In particular, the probability that a point will fall into any (sufficiently nice) set A is equal to the *length* of that set.

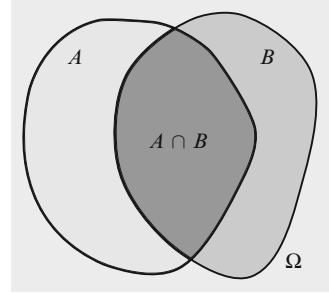
Describing a random experiment by specifying explicitly the sample space and the probability measure is not always straightforward or necessary. Sometimes it is useful to model only certain *observations* on the experiment. This is where *random variables* come into play, and we will discuss these in Chap. 2.

23

1.5 Conditional Probability and Independence

How do probabilities change when we know that some event $B \subset \Omega$ has occurred? Thus, we know that the outcome lies in B . Then A will occur if and only if $A \cap B$ occurs, and the relative chance of A occurring is therefore $\mathbb{P}(A \cap B)/\mathbb{P}(B)$, which is called the *conditional probability* of A given B . The situation is illustrated in Fig. 1.7.

Fig. 1.7 What is the probability that A occurs given that the outcome is known to lie in B ?



Definition 1.4. (Conditional Probability). The **conditional probability** of A given B (with $\mathbb{P}(B) \neq 0$) is defined as:

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} . \quad (1.4)$$

Example 1.7 (Casting Two Dice). We cast two fair dice consecutively. Given that the sum of the dice is 10, what is the probability that one 6 is cast? Let B be the event that the sum is 10:

$$B = \{(4, 6), (5, 5), (6, 4)\} .$$

Let A be the event that one 6 is cast:

$$A = \{(1, 6), \dots, (5, 6), (6, 1), \dots, (6, 5)\} .$$

Then, $A \cap B = \{(4, 6), (6, 4)\}$. And, since for this experiment all elementary events are equally likely, we have

$$\mathbb{P}(A | B) = \frac{2/36}{3/36} = \frac{2}{3} .$$

Example 1.8 (Monty Hall Problem). Consider a quiz in which the final contestant is to choose a prize which is hidden behind one of three curtains (A, B, or C). Suppose without loss of generality that the contestant chooses curtain A. Now the quiz master (Monty) always opens one of the other curtains: if the prize is behind B, Monty opens C; if the prize is behind C, Monty opens B; and if the prize is behind A, Monty opens B or C with equal probability, e.g., by tossing a coin (of course the contestant does not see Monty tossing the coin!) (Fig. 1.8).

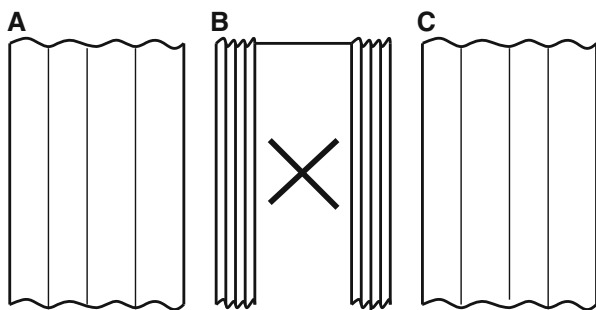
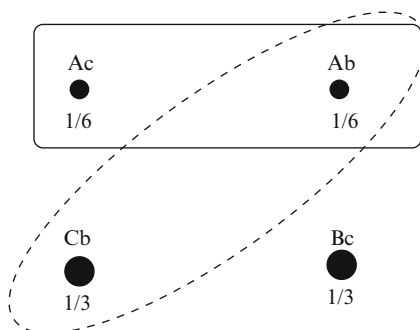


Fig. 1.8 Given that Monty opens curtain (B), should the contestant stay with his/her original choice (A) or switch to the other unopened curtain (C)?

Suppose, again without loss of generality, that Monty opens curtain B. The contestant is now offered the opportunity to switch to curtain C. Should the contestant stay with his/her original choice (A) or switch to the other unopened curtain (C)?

Notice that the sample space here consists of four possible outcomes: Ac , the prize is behind A and Monty opens C; Ab , the prize is behind A and Monty opens B; Bc , the prize is behind B and Monty opens C; and Cb , the prize is behind C and Monty opens B. Let A , B , C be the events that the prize is behind A, B, and C, respectively. Note that $A = \{Ac, Ab\}$, $B = \{Bc\}$, and $C = \{Cb\}$; see Fig. 1.9.

Fig. 1.9 The sample space for the Monty Hall problem



Now, obviously $\mathbb{P}(A) = \mathbb{P}(B) = \mathbb{P}(C)$, and since Ac and Ab are equally likely, we have $\mathbb{P}(\{Ab\}) = \mathbb{P}(\{Ac\}) = 1/6$. Monty opening curtain B means that we have information that event $\{Ab, Cb\}$ has occurred. The probability that the prize is behind A given this event is therefore

$$\mathbb{P}(A \mid B \text{ is opened}) = \frac{\mathbb{P}(\{Ac, Ab\} \cap \{Ab, Cb\})}{\mathbb{P}(\{Ab, Cb\})} = \frac{\mathbb{P}(\{Ab\})}{\mathbb{P}(\{Ab, Cb\})} = \frac{\frac{1}{6}}{\frac{1}{6} + \frac{1}{3}} = \frac{1}{3}.$$

This is what is to be expected: the fact that Monty opens a curtain does not give any extra information that the prize is behind A. Obviously, $\mathbb{P}(B \mid B \text{ is opened}) = 0$. It follows then that $\mathbb{P}(C \mid B \text{ is opened})$ must be $2/3$, since the conditional probabilities must sum up to 1. Indeed,

$$\mathbb{P}(C \mid B \text{ is opened}) = \frac{\mathbb{P}(\{Cb\} \cap \{Ab, Cb\})}{\mathbb{P}(\{Ab, Cb\})} = \frac{\mathbb{P}(\{Cb\})}{\mathbb{P}(\{Ab, Cb\})} = \frac{\frac{1}{3}}{\frac{1}{6} + \frac{1}{3}} = \frac{2}{3}.$$

Hence, given the information that B is opened, it is twice as likely that the prize is behind C than behind A. Thus, the contestant should switch!

1.5.1 Product Rule

By the definition of conditional probability (1.4) we have

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B \mid A).$$

It is not difficult to generalize this to n intersections $A_1 \cap A_2 \cap \cdots \cap A_n$, which we abbreviate as $A_1 A_2 \cdots A_n$. This gives the **product rule** of probability. We leave the proof as an exercise; see Problem 1.11.

20

Theorem 1.3. (Product Rule). Let A_1, \dots, A_n be a sequence of events with $\mathbb{P}(A_1 \cdots A_{n-1}) > 0$. Then,

$$\begin{aligned} \mathbb{P}(A_1 \cdots A_n) &= \\ &\mathbb{P}(A_1) \mathbb{P}(A_2 \mid A_1) \mathbb{P}(A_3 \mid A_1 A_2) \cdots \mathbb{P}(A_n \mid A_1 \cdots A_{n-1}). \end{aligned} \quad (1.5)$$

Example 1.9 (Urn Problem). We draw consecutively three balls from an urn with 5 white and 5 black balls, without putting them back. What is the probability that all drawn balls will be black?

Let A_i be the event that the i th ball is black. We wish to find the probability of $A_1 A_2 A_3$, which by the product rule (1.5) is

$$\mathbb{P}(A_1) \mathbb{P}(A_2 | A_1) \mathbb{P}(A_3 | A_1 A_2) = \frac{5}{10} \frac{4}{9} \frac{3}{8} \approx 0.083 .$$

Example 1.10 (Birthday Problem). What is the probability that in a group of n people all have different birthdays? We can use the product rule. Let A_i be the event that the first i people have different birthdays, $i = 1, 2, \dots$. Note that $\dots \subset A_3 \subset A_2 \subset A_1$. Therefore, $A_n = A_1 \cap A_2 \cap \dots \cap A_n$, and thus by the product rule

$$\mathbb{P}(A_n) = \mathbb{P}(A_1) \mathbb{P}(A_2 | A_1) \mathbb{P}(A_3 | A_2) \cdots \mathbb{P}(A_n | A_{n-1}) .$$

Now $\mathbb{P}(A_k | A_{k-1}) = (365 - k + 1)/365$, because given that the first $k - 1$ people have different birthdays, there are no duplicate birthdays among the first k people if and only if the birthday of the k th person is chosen from the $365 - (k - 1)$ remaining birthdays. Thus, we obtain

$$\mathbb{P}(A_n) = \frac{365}{365} \times \frac{364}{365} \times \frac{363}{365} \times \cdots \times \frac{365 - n + 1}{365}, \quad n \geq 1 . \quad (1.6)$$

A graph of $\mathbb{P}(A_n)$ against n is given in Fig. 1.10. Note that the probability $\mathbb{P}(A_n)$ rapidly decreases to zero. For $n = 23$ the probability of having no duplicate birthdays is already less than $1/2$.

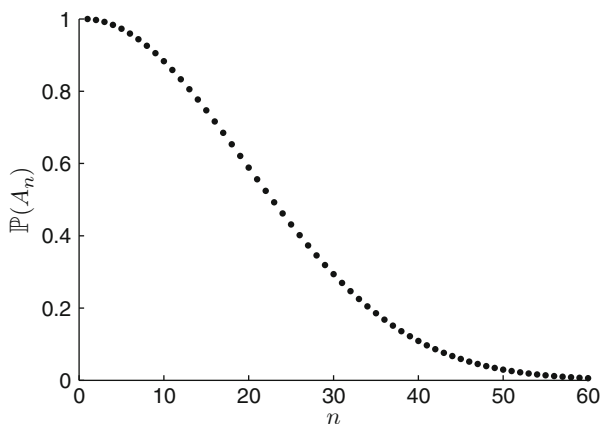


Fig. 1.10 The probability of having no duplicate birthday in a group of n people against n

1.5.2 Law of Total Probability and Bayes' Rule

Suppose that B_1, B_2, \dots, B_n is a **partition** of Ω . That is, B_1, B_2, \dots, B_n are disjoint and their union is Ω ; see Fig. 1.11.

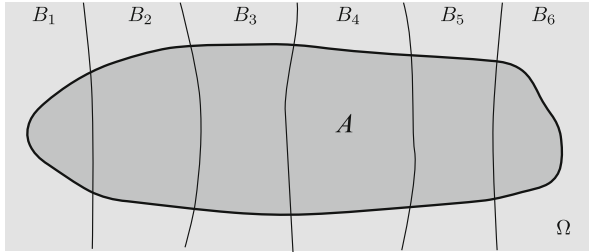


Fig. 1.11 A partition B_1, \dots, B_6 of the sample space Ω . Event A is partitioned into events $A \cap B_1, \dots, A \cap B_6$

A partitioning of the state space can sometimes make it easier to calculate probabilities via the following theorem.

Theorem 1.4. (Law of Total Probability). Let A be an event and let B_1, B_2, \dots, B_n be a partition of Ω . Then,

$$\mathbb{P}(A) = \sum_{i=1}^n \mathbb{P}(A \mid B_i) \mathbb{P}(B_i) . \quad (1.7)$$

Proof. The sum rule gives $\mathbb{P}(A) = \sum_{i=1}^n \mathbb{P}(A \cap B_i)$, and by the product rule we have $\mathbb{P}(A \cap B_i) = \mathbb{P}(A \mid B_i) \mathbb{P}(B_i)$. \square

Combining the law of total probability with the definition of conditional probability gives **Bayes' Rule**:

Theorem 1.5. (Bayes Rule). Let A be an event with $\mathbb{P}(A) > 0$ and let B_1, B_2, \dots, B_n be a partition of Ω . Then,

$$\mathbb{P}(B_j \mid A) = \frac{\mathbb{P}(A \mid B_j) \mathbb{P}(B_j)}{\sum_{i=1}^n \mathbb{P}(A \mid B_i) \mathbb{P}(B_i)} . \quad (1.8)$$

Proof. By definition, $\mathbb{P}(B_j | A) = \mathbb{P}(A \cap B_j) / \mathbb{P}(A) = \mathbb{P}(A | B_j) \mathbb{P}(B_j) / \mathbb{P}(A)$. Now apply the law of total probability to $\mathbb{P}(A)$. \square

Example 1.11 (Quality Control Problem). A company has three factories (1, 2, and 3) that produce the same chip, each producing 15 %, 35 %, and 50 % of the total production. The probability of a faulty chip at factory 1, 2, and 3 is 0.01, 0.05, and 0.02, respectively. Suppose we select randomly a chip from the total production and this chip turns out to be faulty. What is the conditional probability that this chip has been produced in factory 1?

Let B_i denote the event that the chip has been produced in factory i . The $\{B_i\}$ form a partition of Ω . Let A denote the event that the chip is faulty. We are given the information that $\mathbb{P}(B_1) = 0.15$, $\mathbb{P}(B_2) = 0.35$, $\mathbb{P}(B_3) = 0.5$ as well as $\mathbb{P}(A | B_1) = 0.01$, $\mathbb{P}(A | B_2) = 0.05$, $\mathbb{P}(A | B_3) = 0.02$.

We wish to find $\mathbb{P}(B_1 | A)$, which by Bayes' rule is given by

$$\mathbb{P}(B_1 | A) = \frac{0.15 \times 0.01}{0.15 \times 0.01 + 0.35 \times 0.05 + 0.5 \times 0.02} = 0.052 .$$

1.5.3 Independence

Independence is a very important concept in probability and statistics. Loosely speaking it models the *lack of information* between events. We say events A and B are *independent* if the knowledge that B has occurred does not change the probability that A occurs. More precisely, A and B are said to be independent if $\mathbb{P}(A | B) = \mathbb{P}(A)$. Since $\mathbb{P}(A | B) = \mathbb{P}(A \cap B) / \mathbb{P}(B)$, an alternative definition of independence is: A and B are independent if $\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$. This definition covers the case where $B = \emptyset$.

We can extend the definition to arbitrarily many events [compare with the product rule (1.5)]:

Definition 1.5. (Independence). The events A_1, A_2, \dots , are said to be **independent** if for any k and any choice of distinct indices i_1, \dots, i_k ,

$$\mathbb{P}(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_k}) = \mathbb{P}(A_{i_1}) \mathbb{P}(A_{i_2}) \dots \mathbb{P}(A_{i_k}) . \quad (1.9)$$

Remark 1.1. In most cases independence of events is a *model assumption*. That is, \mathbb{P} is chosen such that certain events are independent.

Example 1.12 (Coin Tossing and the Binomial Law). We toss a coin n times. The sample space can be written as the set of binary n -tuples:

$$\Omega = \{(\underbrace{0, \dots, 0}_{n \text{ times}}, \dots, (1, \dots, 1))\}.$$

Here, 0 represents Tails and 1 represents Heads. For example, the outcome $(0, 1, 0, 1, \dots)$ means that the first time Tails is thrown, the second time Heads, the third time Tails, the fourth time Heads, etc.

How should we define \mathbb{P} ? Let A_i denote the event of Heads at the i th throw, $i = 1, \dots, n$. Then, \mathbb{P} should be such that the following holds:

- The events A_1, \dots, A_n should be *independent* under \mathbb{P} .
- $\mathbb{P}(A_i)$ should be the same for all i . Call this known or unknown probability p ($0 \leq p \leq 1$).

These two rules completely specify \mathbb{P} . For example, the probability that the first k throws are Heads and the last $n - k$ are Tails is

$$\begin{aligned} \mathbb{P}(\{(\underbrace{1, 1, \dots, 1}_{k \text{ times}}, \underbrace{0, 0, \dots, 0}_{n-k \text{ times}})\}) &= \mathbb{P}(A_1 \cap \dots \cap A_k \cap A_{k+1}^c \cap \dots \cap A_n^c) \\ &= \mathbb{P}(A_1) \cdots \mathbb{P}(A_k) \mathbb{P}(A_{k+1}^c) \cdots \mathbb{P}(A_n^c) = p^k (1 - p)^{n-k}. \end{aligned}$$

Note that if A_i and A_j are independent, then so are A_i and A_j^c ; see Problem 1.12.

Let B_k be the event that k Heads are thrown in total. The probability of this event is the sum of the probabilities of elementary events $\{(x_1, \dots, x_n)\}$ for which $x_1 + \dots + x_n = k$. Each of these events has probability $p^k (1 - p)^{n-k}$, and there are $\binom{n}{k}$ of these. We thus obtain the **binomial law**:

$$\mathbb{P}(B_k) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, \dots, n. \quad (1.10)$$

Example 1.13 (Geometric Law). There is another important law associated with the coin toss experiment. Let C_k be the event that Heads appears for the first time at the k th toss, $k = 1, 2, \dots$. Then, using the same events $\{A_i\}$ as in the previous example, we can write

$$C_k = A_1^c \cap A_2^c \cap \dots \cap A_{k-1}^c \cap A_k.$$

Using the independence of $A_1^c, \dots, A_{k-1}^c, A_k$, we obtain the **geometric law**:

$$\begin{aligned} \mathbb{P}(C_k) &= \mathbb{P}(A_1^c) \cdots \mathbb{P}(A_{k-1}^c) \mathbb{P}(A_k) \\ &= \underbrace{(1 - p) \cdots (1 - p)}_{k-1 \text{ times}} p = (1 - p)^{k-1} p. \end{aligned}$$

1.6 Problems

1.1. For each of the five random experiments at the beginning of Sect. 1.1 define a convenient sample space.

1.2. Interpret De Morgan's rule (1.2) in terms of an unreliable series system.

1.3. Let $\mathbb{P}(A) = 0.9$ and $\mathbb{P}(B) = 0.8$. Show that $\mathbb{P}(A \cap B) \geq 0.7$.

1.4. Throw two fair dice one after the other.

(a) What is the probability that the second die is 3, given that the sum of the dice is 6?

(b) What is the probability that the first die is 3 and the second is not 3?

1.5. An "expert" wine taster has to try to match 6 glasses of wine to 6 wine labels. Each label can only be chosen once.

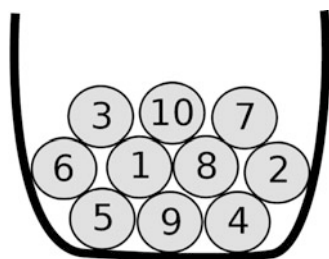
(a) Formulate a sample space Ω for this experiment.

(b) Assuming the wine taster is a complete fraud, define an appropriate probability \mathbb{P} on the sample space.

(c) What is the probability that the wine taster guesses 4 labels correctly, assuming he/she guesses them randomly?

1.6. Many counting problems can be cast into the framework of drawing k balls from an urn with n balls, numbered $1, \dots, n$; see Fig. 1.12.

Fig. 1.12 Draw k balls from an urn with $n = 10$ numbered balls



The drawing can be done in several ways. Firstly, the k balls could be drawn one-by-one or all at the same time. In the first case the **order** in which the balls are drawn can be noted. In the second case we can still assume that the balls are drawn one-by-one, but we do not note the order. Secondly, once a ball is drawn, it can either be put back into the urn or be left out. This is called drawing with and without **replacement**, respectively. There are thus four possible random experiments. Prove that for each of these experiments the total number of possible outcomes is the following:

1. Ordered, with replacement: n^k .

2. Ordered, without replacement: ${}^n P_k = n(n-1) \cdots (n-k+1)$.

3. Unordered, without replacement: ${}^nC_k = \binom{n}{k} = \frac{{}^nP_k}{k!} = \frac{n!}{(n-k)!k!}$.

4. Unordered, with replacement: $\binom{n+k-1}{k}$.

Provide a sample space for each of these experiments. Hint: it is important to use a notation that clearly shows whether the arrangements of numbers are ordered or not. Denote ordered arrangements by *vectors*, e.g., $(1, 1, 2)$, and unordered arrangements by *sets*, e.g., $\{1, 2, 3\}$ or *multisets*, e.g., $\{1, 1, 2\}$.

1.7. Formulate the birthday problem in terms of an urn experiment, as in Problem 1.6, and derive the probability (1.6) by counting.

1.8. Three cards are drawn from a full deck of cards, noting the order. The cards may be numbered from 1 to 52.

- Give the sample space. Is each elementary event equally likely?
- What is the probability that we draw three Aces?
- What is the probability that we draw one Ace, one King, and one Queen (not necessarily in that order)?
- What is the probability that we draw no pictures (no A, K, Q, or J)?

1.9. In a group of 20 people there are three brothers. The group is separated at random into two groups of 10. What is the probability that the brothers are in the same group?

1.10. Two fair dice are thrown.

- Find the probability that both dice show the same face.
- Find the same probability, using the extra information, that the sum of the dice is not greater than 4.

1.11. Prove the product rule (1.5). Hint: first show it for the case of three events:

$$\mathbb{P}(A \cap B \cap C) = \mathbb{P}(A) \mathbb{P}(B | A) \mathbb{P}(C | A \cap B).$$

1.12. If A and B are independent events, then A and B^c are also independent. Prove this.

1.13. Select at random 3 people from a large population. What is the probability that they all have the same birthday?

1.14. In a large population 40 % votes for A and 60 % for B. Suppose we select at random 10 people. What is the probability that in this group exactly 4 people will vote for A?

1.15. A certain AIDS test has a 0.98 probability of giving a positive result when the blood is infected and a 0.07 probability of giving a positive result when the blood is not infected (a so-called false-positive). Suppose 1 % of the population carries the HIV virus.

1. Using the law of total probability, what is the probability that the test is positive for a randomly selected person?
2. What is the probability that a person is indeed infected, *given* that the test yields a positive result?

1.16. A box has three identical-looking coins. However, the probability of success (Heads) is different for each coin: coin 1 is fair, coin 2 has a success probability of 0.4, and coin 3 has a success probability of 0.6. We pick one coin at random and throw it 100 times. Suppose 43 Heads come up. Using this information, assess the probability that coin 1, 2, or 3 was chosen.

1.17. In a binary communication channel, 0s and 1s are transmitted with equal probability. The probability that a 0 is correctly received (as a 0) is 0.95. The probability that a 1 is correctly received (as a 1) is 0.99. Suppose we receive a 0, what is the probability that, in fact, a 1 was sent?

1.18. A fair coin is tossed 20 times.

1. What is the probability of exactly 10 Heads?
2. What is the probability of 15 or more Heads?

1.19. Two fair dice are cast (at the same time) until their sum is 12.

1. What is the probability that we have to wait exactly 10 tosses?
2. What is the probability that we do not have to wait more than 100 tosses?

1.20. Independently throw 10 balls into one of three boxes, numbered 1, 2, and 3, with probabilities $1/4$, $1/2$, and $1/4$, respectively.

1. What is the probability that box 1 has 2 balls, box 2 has 5 balls, and box 3 has 3 balls?
2. What is the probability that box 1 remains empty?



1.21. Implement a MATLAB program that performs 100 tosses with a fair die. Hint: use the `rand` and `ceil` functions, where `ceil(x)` returns the smallest integer larger than or equal to x .



1.22. For each of the four urn experiments in Problem 1.6 implement a MATLAB program that simulates the experiment. Hint: in addition to the functions `rand` and `ceil`, you may wish to use the `sort` function.



1.23. Verify your answers for Problem 1.20 with a computer simulation, where the experiment is repeated many times.

Chapter 2

Random Variables and Probability Distributions

Specifying a model for a random experiment via a complete description of the sample space Ω and probability measure \mathbb{P} may not always be necessary or convenient. In practice we are only interested in certain *numerical measurements* pertaining to the experiment. Such random measurements can be included into the model via the notion of a *random variable*.

2.1 Random Variables

Definition 2.1. (Random Variable). A **random variable** is a *function* from the sample space Ω to \mathbb{R} .

Example 2.1 (Sum of Two Dice). We throw two fair dice and note the sum of their face values. If we throw the dice consecutively and observe both throws, the sample space is $\Omega = \{(1, 1), \dots, (6, 6)\}$. The function X defined by $X(i, j) = i + j$ is a random variable which maps the outcome (i, j) to the sum $i + j$, as depicted in Fig. 2.1.

Note that five outcomes in the sample space are mapped to 8. A natural notation for the corresponding set of outcomes is $\{X = 8\}$. Since all outcomes in Ω are equally likely, we have

$$\mathbb{P}(\{X = 8\}) = \frac{5}{36}.$$

This notation is very suggestive and convenient. From a nonmathematical viewpoint we can interpret X as a “random” variable, that is, a variable that can take several values with certain probabilities. In particular, it is not difficult to check that

$$\mathbb{P}(\{X = x\}) = \frac{6 - |7 - x|}{36}, \quad x = 2, \dots, 12.$$

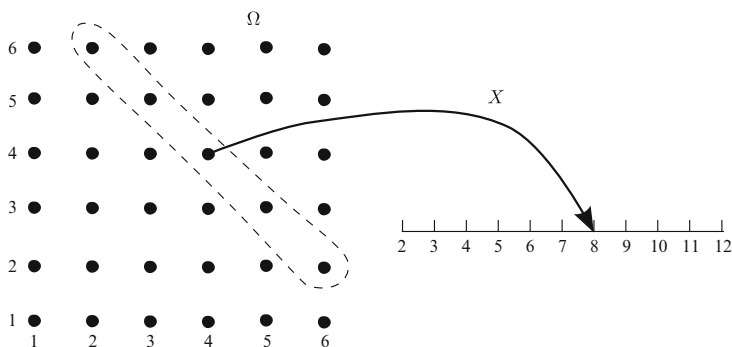


Fig. 2.1 Random variable X represents the sum of two dice

Although random variables are, mathematically speaking, *functions*, it is often convenient to view them as observations of a random experiment that has not yet taken place. In other words, a random variable is considered as a measurement that becomes available *tomorrow*, while all the thinking about the measurement can be carried out *today*. For example, we can specify today exactly the probabilities pertaining to the random variables.

We often denote random variables with *capital* letters from the last part of the alphabet, e.g., X, X_1, X_2, \dots, Y, Z . Random variables allow us to use natural and intuitive notations for certain events, such as $\{X = 10\}$, $\{X > 1000\}$, and $\{\max(X, Y) \leq Z\}$.

18 Example 2.2 (Coin Tossing). In Example 1.12 we constructed a probability model for the random experiment where a biased coin is tossed n times. Suppose we are not interested in a specific outcome but only in the total number of Heads, X , say. In particular, we would like to know the probability that X takes certain values between 0 and n . Example 1.12 suggests that

$$\mathbb{P}(X = k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, \dots, n, \quad (2.1)$$

providing all the information about X that we could possibly wish to know. To justify (2.1) mathematically, we can reason as in Example 2.1. First, define X as the function that assigns to each outcome $\omega = (x_1, \dots, x_n)$ the number $x_1 + \dots + x_n$. Thus, X is a random variable in mathematical terms, that is, a function. Second, the event B_k that there are exactly k Heads in n throws can be written as

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

If we write this as $\{X = k\}$, and further abbreviate $\mathbb{P}(\{X = k\})$ to $\mathbb{P}(X = k)$, then we obtain (2.1) directly from (1.10).

We give some more examples of random variables without specifying the sample space:

1. The number of defective transistors out of 100 inspected ones.
2. The number of bugs in a computer program.
3. The amount of rain in a certain location in June.
4. The amount of time needed for an operation.

The set of all possible values that a random variable X can take is called the **range** of X . We further distinguish between discrete and continuous random variables:

- **Discrete** random variables can only take *countably many* values.
- **Continuous** random variables can take a continuous range of values, for example, any value on the positive real line \mathbb{R}_+ .

2.2 Probability Distribution

Let X be a random variable. We would like to designate the probabilities of events such as $\{X = x\}$ and $\{a \leq X \leq b\}$. If we can specify all probabilities involving X , we say that we have determined the **probability distribution** of X . One way to specify the probability distribution is to give the probabilities of all events of the form $\{X \leq x\}$, $x \in \mathbb{R}$. This leads to the following definition.

Definition 2.2. (Cumulative Distribution Function). The **cumulative distribution function** (cdf) of a random variable X is the function $F : \mathbb{R} \rightarrow [0, 1]$ defined by

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

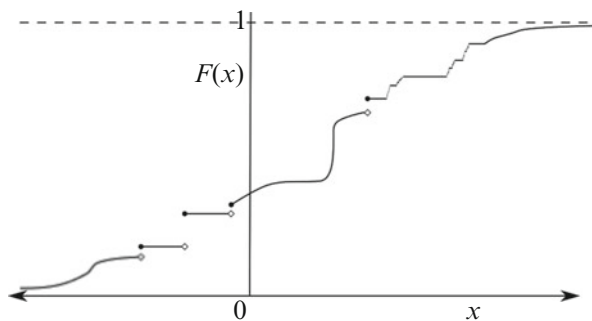


Fig. 2.2 A cumulative distribution function (cdf)

Note that we have used $\mathbb{P}(X \leq x)$ as a shorthand notation for $\mathbb{P}(\{X \leq x\})$. From now on we will use this type of abbreviation throughout the book. In Fig. 2.2 the graph of a general cdf is depicted.

Theorem 2.1. (Properties of Cdf). Let F be the cdf of a random variable X . Then,

1. F is bounded between 0 and 1: $0 \leq F(x) \leq 1$,
2. F is increasing: if $x \leq y$, then $F(x) \leq F(y)$,
3. F is right-continuous: $\lim_{h \downarrow 0} F(x+h) = F(x)$.

Proof.

- 9 1. Let $A = \{X \leq x\}$. By Rule 1 in Definition 1.3, $0 \leq \mathbb{P}(A) \leq 1$.
- 10 2. Suppose $x \leq y$. Define $A = \{X \leq x\}$ and $B = \{X \leq y\}$. Then, $A \subset B$ and, by Theorem 1.2, $\mathbb{P}(A) \leq \mathbb{P}(B)$.
3. Take any sequence x_1, x_2, \dots decreasing to x . We have to show that $\lim_{n \rightarrow \infty} \mathbb{P}(X \leq x_n) = \mathbb{P}(X \leq x)$ or, equivalently, $\lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \mathbb{P}(A)$, where $A_n = \{X > x_n\}$ and $A = \{X > x\}$. Let $B_n = \{x_{n-1} \geq X > x_n\}$, $n = 1, 2, \dots$, with x_0 defined as ∞ . Then, $A_n = \cup_{i=1}^n B_i$ and $A = \cup_{i=1}^{\infty} B_i$. Since the $\{B_i\}$ are disjoint, we have by the sum rule:

$$\mathbb{P}(A) = \sum_{i=1}^{\infty} \mathbb{P}(B_i) \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{P}(B_i) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n) ,$$

as had to be shown. □

Conversely, any function F with the above properties can be used to specify the distribution of a random variable X .

If X has cdf F , then the probability that X takes a value in the interval $(a, b]$ (excluding a , including b) is given by

$$\mathbb{P}(a < X \leq b) = F(b) - F(a) .$$

To see this, note that $\mathbb{P}(X \leq b) = \mathbb{P}(\{X \leq a\} \cup \{a < X \leq b\})$, where the events $\{X \leq a\}$ and $\{a < X \leq b\}$ are disjoint. Thus, by the sum rule, $F(b) = F(a) + \mathbb{P}(a < X \leq b)$, which leads to the result above. Note however that

$$\begin{aligned} \mathbb{P}(a \leq X \leq b) &= F(b) - F(a) + \mathbb{P}(X = a) \\ &= F(b) - F(a) + F(a) - F(a-) \\ &= F(b) - F(a-) , \end{aligned}$$

where $F(a-)$ denotes the limit from below: $\lim_{x \uparrow a} F(x)$.

2.2.1 Discrete Distributions

Definition 2.3. (Discrete Distribution). A random variable X is said to have a **discrete distribution** if $\mathbb{P}(X = x_i) > 0, i = 1, 2, \dots$ for some finite or countable set of values x_1, x_2, \dots , such that $\sum_i \mathbb{P}(X = x_i) = 1$. The **discrete probability density function (pdf)** of X is the function f defined by $f(x) = \mathbb{P}(X = x)$.

We sometimes write f_X instead of f to stress that the discrete probability density function refers to the discrete random variable X . The easiest way to specify the distribution of a discrete random variable is to specify its pdf. Indeed, by the sum rule, if we know $f(x)$ for all x , then we can calculate all possible probabilities involving X . Namely,

9

$$\mathbb{P}(X \in B) = \sum_{x \in B} f(x) \tag{2.2}$$

for any subset B in the range of X , as illustrated in Fig. 2.3.

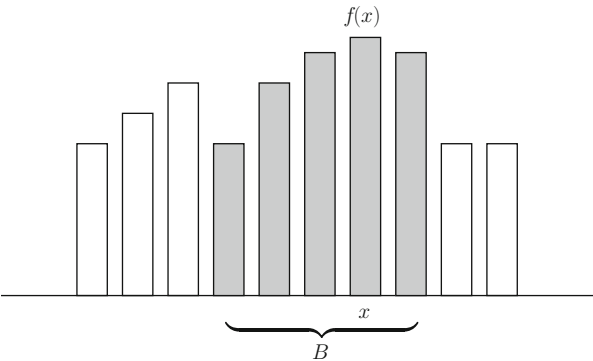


Fig. 2.3 Discrete probability density function

Example 2.3 (Sum of Two Dice, Continued). Toss two fair dice and let X be the sum of their face values. The discrete pdf is given in Table 2.1, which follows directly from Example 2.1.

Table 2.1 Discrete pdf of the sum of two fair dice

x	2	3	4	5	6	7	8	9	10	11	12
$f(x)$	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$

2.2.2 Continuous Distributions

Definition 2.4. (Continuous Distribution). A random variable X with cdf F is said to have a **continuous distribution** if there exists a positive function f with *total integral 1* such that for all $a < b$,

$$\mathbb{P}(a < X \leq b) = F(b) - F(a) = \int_a^b f(u) du . \quad (2.3)$$

Function f is called the **probability density function (pdf)** of X .

Remark 2.1. Note that we use the *same* notation f for both the discrete and the continuous pdf, to stress the similarities between the discrete and continuous case. We will even drop the qualifier “discrete” or “continuous” when it is clear from the context with which case we are dealing. Henceforth we will use the notation $X \sim f$ and $X \sim F$ to indicate that X is distributed according to pdf f or cdf F .

In analogy to the discrete case (2.2), once we know the pdf, we can calculate any probability of interest by means of integration:

$$\mathbb{P}(X \in B) = \int_B f(x) dx , \quad (2.4)$$

as illustrated in Fig. 2.4.

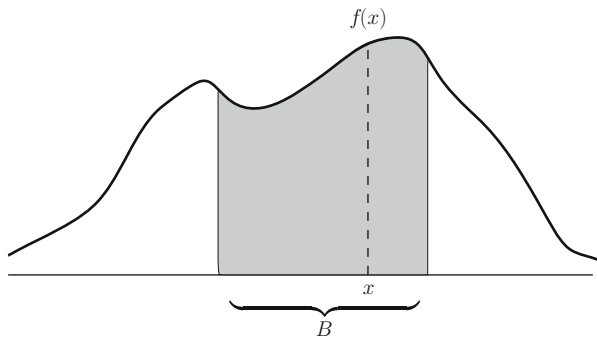


Fig. 2.4 Probability density function (pdf)

Suppose that f and F are the pdf and cdf of a continuous random variable X , as in Definition 2.4. Then F is simply a *primitive* (also called antiderivative) of f :

$$F(x) = \mathbb{P}(X \leq x) = \int_{-\infty}^x f(u) du .$$

Conversely, f is the *derivative* of the cdf F :

$$f(x) = \frac{d}{dx} F(x) = F'(x) .$$

It is important to understand that in the continuous case $f(x)$ is not equal to the probability $\mathbb{P}(X = x)$, because the latter is 0 for all x . Instead, we interpret $f(x)$ as the *density* of the probability distribution at x , in the sense that for any small h ,

$$\mathbb{P}(x \leq X \leq x + h) = \int_x^{x+h} f(u) du \approx h f(x) . \quad (2.5)$$

Note that $\mathbb{P}(x \leq X \leq x + h)$ is equal to $\mathbb{P}(x < X \leq x + h)$ in this case.

Example 2.4 (Random Point in an Interval). Draw a random number X from the interval of real numbers $[0, 2]$, where each number is equally likely to be drawn. What are the pdf f and cdf F of X ? Using the same reasoning as in Example 1.6, we see that

 11

$$\mathbb{P}(X \leq x) = F(x) = \begin{cases} 0 & \text{if } x < 0, \\ x/2 & \text{if } 0 \leq x \leq 2, \\ 1 & \text{if } x > 2. \end{cases}$$

By differentiating F we find

$$f(x) = \begin{cases} 1/2 & \text{if } 0 \leq x \leq 2, \\ 0 & \text{otherwise.} \end{cases}$$

Note that this density is *constant* on the interval $[0, 2]$ (and zero elsewhere), reflecting the fact that each point in $[0, 2]$ is equally likely to be drawn.

2.3 Expectation

Although all probability information about a random variable is contained in its cdf or pdf, it is often useful to consider various numerical characteristics of a random variable. One such number is the *expectation* of a random variable, which is a “weighted average” of the values that X can take. Here is a more precise definition.

Definition 2.5. (Expectation of a Discrete Random Variable). Let X be a *discrete* random variable with pdf f . The **expectation** (or expected value) of X , denoted as $\mathbb{E}X$, is defined as

$$\mathbb{E}X = \sum_x x \mathbb{P}(X = x) = \sum_x x f(x) . \quad (2.6)$$

The expectation of X is sometimes written as μ_X . It is assumed that the sum in (2.6) is well defined—possibly ∞ or $-\infty$. One way to interpret the expectation is as a *long-run average payout*. Suppose in a game of dice the payout X (dollars) is the largest of the face values of two dice. To play the game a fee of d dollars must be paid. What would be a fair amount for d ? The answer is

$$\begin{aligned} d = \mathbb{E}X &= 1 \times \mathbb{P}(X = 1) + 2 \times \mathbb{P}(X = 2) + \cdots + 6 \times \mathbb{P}(X = 6) \\ &= 1 \times \frac{1}{36} + 2 \times \frac{3}{36} + 3 \times \frac{5}{36} + 4 \times \frac{7}{36} + 5 \times \frac{9}{36} + 6 \times \frac{11}{36} = \frac{161}{36} \approx 4.47. \end{aligned}$$

Namely, if the game is played many times, the long-run fraction of tosses in which the maximum face value is 1, 2, ..., 6, is $\frac{1}{36}, \frac{3}{36}, \dots, \frac{11}{36}$, respectively. Hence, the long-run average payout of the game is the weighted sum of 1, 2, ..., 6, where the weights are the long-run fractions (probabilities). The game is “fair” if the long-run average profit $\mathbb{E}X - d$ is zero.

The expectation can also be interpreted as a *center of mass*. Imagine that point masses with weights p_1, p_2, \dots, p_n are placed at positions x_1, x_2, \dots, x_n on the real line; see Fig. 2.5.

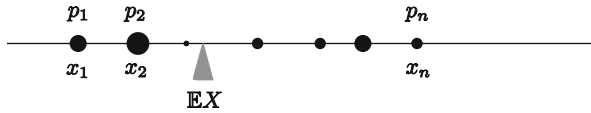


Fig. 2.5 The expectation as a center of mass

The center of mass—the place where the weights are balanced—is

$$\text{center of mass} = x_1 p_1 + \cdots + x_n p_n,$$

which is exactly the expectation of the discrete variable X that takes values x_1, \dots, x_n with probabilities p_1, \dots, p_n . An obvious consequence of this interpretation is that for a *symmetric* pdf the expectation is equal to the symmetry point (provided that the expectation exists). In particular, suppose that $f(c + y) = f(c - y)$ for all y . Then,

$$\begin{aligned} \mathbb{E}X &= c f(c) + \sum_{x>c} x f(x) + \sum_{x<c} x f(x) \\ &= c f(c) + \sum_{y>0} (c + y) f(c + y) + \sum_{y>0} (c - y) f(c - y) \\ &= c f(c) + \sum_{y>0} c f(c + y) + c \sum_{y>0} f(c - y) = c \sum_x f(x) = c. \end{aligned}$$

For continuous random variables we can define the expectation in a similar way, replacing the sum with an integral.

Definition 2.6. (Expectation of a Continuous Random Variable). Let X be a *continuous* random variable with pdf f . The **expectation** (or expected value) of X , denoted as $\mathbb{E}X$, is defined as

$$\mathbb{E}X = \int_{-\infty}^{\infty} x f(x) dx . \quad (2.7)$$

If X is a random variable, then a function of X , such as X^2 or $\sin(X)$, is also a random variable. The following theorem simply states that the expected value of a function of X is the weighted average of the values that this function can take.

Theorem 2.2. (Expectation of a Function of a Random Variable). If X is *discrete* with pdf f , then for any real-valued function g

$$\mathbb{E} g(X) = \sum_x g(x) f(x) .$$

Similarly, if X is *continuous* with pdf f , then

$$\mathbb{E} g(X) = \int_{-\infty}^{\infty} g(x) f(x) dx .$$

Proof. The proof is given for the discrete case only, as the continuous case can be proven in a similar way. Let $Y = g(X)$, where X is a discrete random variable with pdf f_X and g is a function. Let f_Y be the (discrete) pdf of the random variable Y . It can be expressed in terms of f_X in the following way:

$$f_Y(y) = \mathbb{P}(Y = y) = \mathbb{P}(g(X) = y) = \sum_{x: g(x)=y} \mathbb{P}(X = x) = \sum_{x: g(x)=y} f_X(x) .$$

Thus, the expectation of Y is

$$\begin{aligned} \mathbb{E}Y &= \sum_y y f_Y(y) = \sum_y y \sum_{x: g(x)=y} f_X(x) = \sum_y \sum_{x: g(x)=y} y f_X(x) \\ &= \sum_x g(x) f_X(x) . \end{aligned}$$

□

Example 2.5 (Die Experiment and Expectation). Find $\mathbb{E}X^2$ if X is the outcome of the toss of a fair die. We have

$$\mathbb{E}X^2 = 1^2 \times \frac{1}{6} + 2^2 \times \frac{1}{6} + 3^2 \times \frac{1}{6} + \cdots + 6^2 \times \frac{1}{6} = \frac{91}{6}.$$

An important consequence of Theorem 2.2 is that the expectation is “linear”.

Theorem 2.3. (Properties of the Expectation). For any real numbers a and b , and functions g and h ,

1. $\mathbb{E}[aX + b] = a\mathbb{E}X + b$,
2. $\mathbb{E}[g(X) + h(X)] = \mathbb{E}g(X) + \mathbb{E}h(X)$.

Proof. Suppose X has pdf f . The first statement follows (in the discrete case) from

$$\mathbb{E}(aX + b) = \sum_x (ax + b)f(x) = a \sum_x x f(x) + b \sum_x f(x) = a\mathbb{E}X + b.$$

Similarly, the second statement follows from

$$\begin{aligned} \mathbb{E}(g(X) + h(X)) &= \sum_x (g(x) + h(x))f(x) = \sum_x g(x)f(x) + \sum_x h(x)f(x) \\ &= \mathbb{E}g(X) + \mathbb{E}h(X). \end{aligned}$$

The continuous case is proven analogously, simply by replacing sums with integrals. \square

Another useful numerical characteristic of the distribution of X is the *variance* of X . This number, sometimes written as σ_X^2 , measures the *spread* or dispersion of the distribution of X .

Definition 2.7. (Variance and Standard Deviation). The **variance** of a random variable X , denoted as $\text{Var}(X)$, is defined as

$$\text{Var}(X) = \mathbb{E}(X - \mathbb{E}X)^2. \quad (2.8)$$

The square root of the variance is called the **standard deviation**. The number $\mathbb{E}X^r$ is called the r th **moment** of X .

Theorem 2.4. (Properties of the Variance). For any random variable X the following properties hold for the variance:

1. $\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2$.
2. $\text{Var}(a + bX) = b^2 \text{Var}(X)$.

Proof. Write $\mathbb{E}X = \mu$, so that $\text{Var}(X) = \mathbb{E}(X - \mu)^2 = \mathbb{E}(X^2 - 2\mu X + \mu^2)$. By the linearity of the expectation, the last expectation is equal to the sum $\mathbb{E}X^2 - 2\mu \mathbb{E}X + \mu^2 = \mathbb{E}X^2 - \mu^2$, which proves the first statement. To prove the second statement, note that the expectation of $a + bX$ is equal to $a + b\mu$. Consequently,

$$\text{Var}(a + bX) = \mathbb{E}(a + bX - (a + b\mu))^2 = \mathbb{E}(b^2(X - \mu)^2) = b^2 \text{Var}(X).$$

□

Note that Property 1 in Theorem 2.4 implies that $\mathbb{E}X^2 \geq (\mathbb{E}X)^2$, because $\text{Var}(X) \geq 0$. This is a special case of a much more general result, regarding the expectation of convex functions. A real-valued function $h(x)$ is said to be **convex** if for each x_0 , there exist constants a and b such that (1) $h(x) \geq ax + b$ for all x and (2) $h(x_0) = ax_0 + b$. Examples are the functions $x \mapsto x^2$, $x \mapsto e^x$, and $x \mapsto -\ln x$.

Theorem 2.5. (Jensen's Inequality). Let $h(x)$ be a convex function and X a random variable. Then,

$$\mathbb{E}h(X) \geq h(\mathbb{E}X). \quad (2.9)$$

Proof. Let $x_0 = \mathbb{E}X$. Because h is convex, there exist constants a and b such that $h(X) \geq aX + b$ and $h(x_0) = ax_0 + b$. Hence, $\mathbb{E}h(X) \geq \mathbb{E}(aX + b) = ax_0 + b = h(x_0) = h(\mathbb{E}X)$. □

2.4 Transforms

Many probability calculations—such as the evaluation of expectations and variances—are facilitated by the use of *transforms*. We discuss here a number of such transforms.

Definition 2.8. (Probability Generating Function). Let X be a *nonnegative* and *integer-valued* random variable with discrete pdf f . The **probability generating function** (PGF) of X is the function G defined by

$$G(z) = \mathbb{E} z^X = \sum_{x=0}^{\infty} z^x f(x), \quad |z| < R,$$

where $R \geq 1$ is the **radius of convergence**.

Example 2.6 (Poisson Distribution). Let X have a discrete pdf f given by

$$f(x) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

X is said to have a **Poisson distribution**. The PGF of X is given by


$$\begin{aligned} G(z) &= \sum_{x=0}^{\infty} z^x e^{-\lambda} \frac{\lambda^x}{x!} \\ &= e^{-\lambda} \sum_{x=0}^{\infty} \frac{(z\lambda)^x}{x!} \\ &= e^{-\lambda} e^{z\lambda} = e^{-\lambda(1-z)}. \end{aligned}$$

As this is finite for every z , the radius of convergence is here $R = \infty$.

Theorem 2.6. (Derivatives of a PGF). The k th derivative of a PGF $\mathbb{E} z^X$ can be obtained by *differentiation under the expectation sign*:

$$\begin{aligned} \frac{d^k}{dz^k} \mathbb{E} z^X &= \mathbb{E} \frac{d^k}{dz^k} z^X \\ &= \mathbb{E} [X(X-1) \cdots (X-k+1) z^{X-k}] \quad \text{for } |z| < R, \end{aligned}$$

where $R \geq 1$ is the radius of convergence of the PGF.

 **369** *Proof.* The proof is deferred to Appendix B.2. □

Let $G(z)$ be the PGF of a random variable X . Thus, $G(z) = z^0 \mathbb{P}(X = 0) + z^1 \mathbb{P}(X = 1) + z^2 \mathbb{P}(X = 2) + \cdots$. Substituting $z = 0$ gives $G(0) = \mathbb{P}(X = 0)$. By Theorem 2.6 the derivative of G is

$$G'(z) = \mathbb{P}(X = 1) + 2z \mathbb{P}(X = 2) + 3z^2 \mathbb{P}(X = 3) + \cdots$$

In particular, $G'(0) = \mathbb{P}(X = 1)$. By differentiating $G'(z)$, we see that the second derivative of G at 0 is $G''(0) = 2\mathbb{P}(X = 2)$. Repeating this procedure gives the following corollary to Theorem 2.6.

Corollary 2.1. (Probabilities from PGFs). Let X be a nonnegative integer-valued random variable with PGF $G(z)$. Then,

$$\mathbb{P}(X = k) = \frac{1}{k!} \frac{d^k}{dz^k} G(0) .$$

The PGF thus uniquely determines the discrete pdf. Another consequence of Theorem 2.6 is that expectations, variances, and moments can be easily found from the PGF.

Corollary 2.2. (Moments from PGFs). Let X be a nonnegative integer-valued random variable with PGF $G(z)$ and k th derivative $G^{(k)}(z)$. Then,

$$\lim_{\substack{z \rightarrow 1 \\ |z| < 1}} \frac{d^k}{dz^k} G(z) = \mathbb{E}[X(X-1)\cdots(X-k+1)] . \quad (2.10)$$

In particular, if the expectation and variance of X are finite, then $\mathbb{E}X = G'(1)$ and $\text{Var}(X) = G''(1) + G'(1) - (G'(1))^2$.

Proof. The proof is deferred to Appendix B.2. □

 369

Definition 2.9. (Moment Generating Function). The **moment generating function** (MGF) of a random variable X is the function $M : \mathbb{R} \rightarrow [0, \infty]$ given by

$$M(s) = \mathbb{E} e^{sX} .$$

In particular, for a discrete random variable with pdf f ,

$$M(s) = \sum_x e^{sx} f(x) ,$$

and for a continuous random variable with pdf f ,

$$M(s) = \int_{-\infty}^{\infty} e^{sx} f(x) dx .$$

Note that $M(s)$ can be infinite for certain values of s . We sometimes write M_X to stress the role of X .

Similar to the PGF, the MGF has the **uniqueness property**: two MGFs are the same if and only if their corresponding cdfs are the same. In addition, the integer moments of X can be computed from the derivatives of M , as summarized in the next theorem. The proof is similar to that of Theorem 2.6 and Corollary 2.2 and is given in Appendix B.3.

370

Theorem 2.7. (Moments from MGFs). If the MGF is finite in an open interval containing 0, then all moments $\mathbb{E}X^n$, $n = 0, 1, \dots$ are finite and satisfy

$$\mathbb{E}X^n = M^{(n)}(0) ,$$

where $M^{(n)}(0)$ is the n th derivative of M evaluated at 0.

Note that under the conditions of Theorem 2.7, the variance of X can be obtained from the MGF as

$$\text{Var}(X) = M''(0) - (M'(0))^2 .$$

2.5 Common Discrete Distributions

In this section we give a number of common discrete distributions and list some of their properties. Note that the discrete pdf of each of these distributions, denoted f , depends on one or more parameters; so in fact we are dealing with *families* of distributions.

2.5.1 Bernoulli Distribution

Definition 2.10. (Bernoulli Distribution). A random variable X is said to have a **Bernoulli** distribution with success probability p if X can only assume the values 0 and 1, with probabilities

$$f(0) = \mathbb{P}(X = 0) = 1 - p \quad \text{and} \quad f(1) = \mathbb{P}(X = 1) = p .$$

We write $X \sim \text{Ber}(p)$.

The Bernoulli distribution is the most fundamental of all probability distributions. It models a single coin toss experiment. Three important properties of the Bernoulli are summarized in the following theorem.

Theorem 2.8. (Properties of the Bernoulli Distribution). Let $X \sim \text{Ber}(p)$. Then,

1. $\mathbb{E}X = p$,
2. $\text{Var}(X) = p(1 - p)$,
3. the PGF is $G(z) = 1 - p + zp$.

Proof. The expectation and the variance of X can be obtained by direct computation:

$$\mathbb{E}X = 0 \times \mathbb{P}(X = 0) + 1 \times \mathbb{P}(X = 1) = 0 \times (1 - p) + 1 \times p = p$$

and

$$\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \mathbb{E}X - (\mathbb{E}X)^2 = p - p^2 = p(1 - p),$$

where we have used the fact that in this case $X^2 = X$. Finally, the PGF is given by $G(z) = z^0(1 - p) + z^1p = 1 - p + zp$. \square

2.5.2 Binomial Distribution

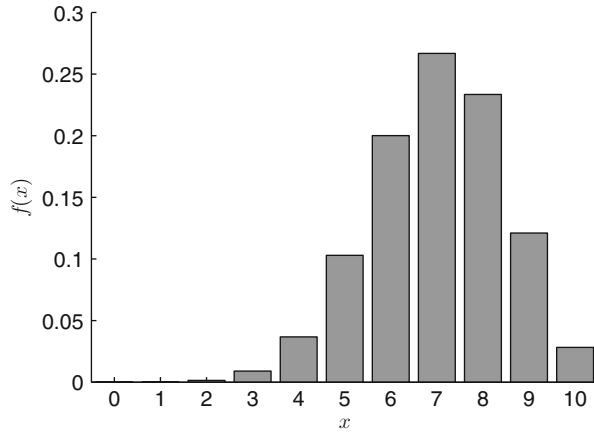
Definition 2.11. (Binomial Distribution). A random variable X is said to have a **binomial** distribution with parameters n and p if X has pdf

$$f(x) = \mathbb{P}(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}, \quad x = 0, 1, \dots, n. \quad (2.11)$$

We write $X \sim \text{Bin}(n, p)$.

From Example 2.2 we see that X can be interpreted as the total number of Heads in n successive coin flip experiments, with probability of Heads equal to p . An example of the graph of the pdf is given in Fig. 2.6. Theorem 2.9 lists some important properties of the binomial distribution.

Fig. 2.6 The pdf of the $\text{Bin}(10, 0.7)$ distribution



Theorem 2.9. (Properties of the Binomial Distribution). Let $X \sim \text{Bin}(n, p)$. Then,

1. $\mathbb{E}X = np$,
2. $\text{Var}(X) = np(1 - p)$,
3. the PGF is $G(z) = (1 - p + zp)^n$.

Proof. Using Newton's binomial formula:

$$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k},$$

we see that

$$G(z) = \sum_{k=0}^n z^k \binom{n}{k} p^k (1 - p)^{n-k} = \sum_{k=0}^n \binom{n}{k} (zp)^k (1 - p)^{n-k} = (1 - p + zp)^n.$$

35 From Corollary 2.2 we obtain the expectation and variance via $G'(1) = np$ and $G''(1) + G'(1) - (G'(1))^2 = (n - 1)np^2 + np - n^2p^2 = np(1 - p)$. \square

2.5.3 Geometric Distribution

Definition 2.12. (Geometric Distribution). A random variable X is said to have a **geometric** distribution with parameter p if X has pdf

$$f(x) = \mathbb{P}(X = x) = (1 - p)^{x-1} p, \quad x = 1, 2, 3, \dots \quad (2.12)$$

We write $X \sim \text{Geom}(p)$.

From Example 1.13 we see that X can be interpreted as the number of tosses needed until the first Heads occurs in a sequence of coin tosses, with the probability of Heads equal to p . An example of the graph of the pdf is given in Fig. 2.7. Theorem 2.10 summarizes some properties of the geometric distribution.

18

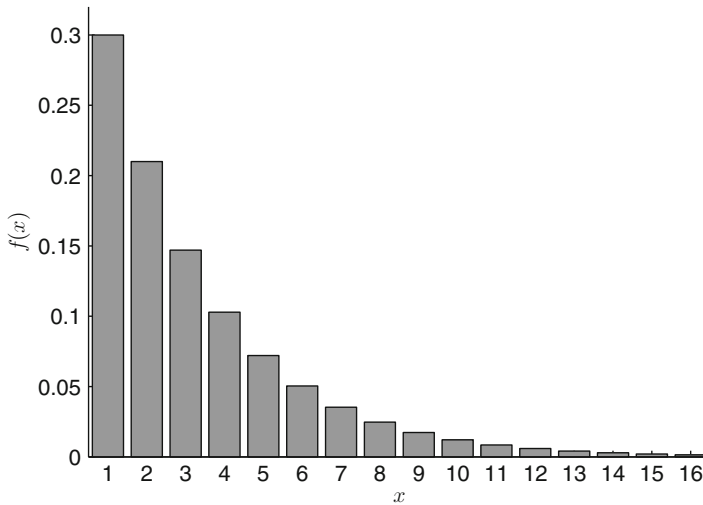


Fig. 2.7 The pdf of the Geom(0.3) distribution

Theorem 2.10. (Properties of the Geometric Distribution). Let $X \sim \text{Geom}(p)$. Then,

1. $\mathbb{E}X = 1/p$,
2. $\text{Var}(X) = (1-p)/p^2$,
3. the PGF is

$$G(z) = \frac{zp}{1 - z(1-p)}, \quad |z| < \frac{1}{1-p}. \quad (2.13)$$

Proof. The PGF of X follows from

$$G(z) = \sum_{x=1}^{\infty} z^x p(1-p)^{x-1} = zp \sum_{k=0}^{\infty} (z(1-p))^k = \frac{zp}{1 - z(1-p)},$$

using the well-known result for *geometric sums*: $1 + a + a^2 + \cdots = (1-a)^{-1}$, for $|a| < 1$. By Corollary 2.2 the expectation is therefore

35

$$\mathbb{E}X = G'(1) = \frac{1}{p}.$$


By differentiating the PGF twice we find the variance:

$$\text{Var}(X) = G''(1) + G'(1) - (G'(1))^2 = \frac{2(1-p)}{p^2} + \frac{1}{p} - \frac{1}{p^2} = \frac{1-p}{p^2}. \quad \square$$

One property of the geometric distribution that deserves extra attention is the **memoryless property**. Consider again the coin toss experiment. Suppose we have tossed the coin k times without a success (Heads). What is the probability that we need more than x additional tosses before getting a success? The answer is, obviously, the same as the probability that we require more than x tosses if we start from scratch, that is, $\mathbb{P}(X > x) = (1-p)^x$, irrespective of k . The fact that we have already had k failures does not make the event of getting a success in the next trial(s) any more likely. In other words, the coin does not have a memory of what happened—hence the name memoryless property.

Theorem 2.11. (Memoryless Property). Let $X \sim \text{Geom}(p)$. Then for any $x, k = 1, 2, \dots$,

$$\mathbb{P}(X > k + x \mid X > k) = \mathbb{P}(X > x).$$

 **12** *Proof.* By the definition of conditional probability,

$$\mathbb{P}(X > k + x \mid X > k) = \frac{\mathbb{P}(\{X > k + x\} \cap \{X > k\})}{\mathbb{P}(X > k)}.$$

The event $\{X > k + x\}$ is a subset of $\{X > k\}$; hence their intersection is $\{X > k + x\}$. Moreover, the probabilities of the events $\{X > k + x\}$ and $\{X > k\}$ are $(1-p)^{k+x}$ and $(1-p)^k$, respectively. Therefore,

$$\mathbb{P}(X > k + x \mid X > k) = \frac{(1-p)^{k+x}}{(1-p)^k} = (1-p)^x = \mathbb{P}(X > x),$$

as required. \square

2.5.4 Poisson Distribution

Definition 2.13. (Poisson Distribution). A random variable X is said to have a **Poisson** distribution with parameter $\lambda > 0$ if X has pdf

(continued)

(continued)

$$f(x) = \mathbb{P}(X = x) = \frac{\lambda^x}{x!} e^{-\lambda}, \quad x = 0, 1, 2, \dots \quad (2.14)$$

We write $X \sim \text{Poi}(\lambda)$.

The Poisson distribution may be viewed as the limit of the $\text{Bin}(n, \lambda/n)$ distribution. Namely, if $X_n \sim \text{Bin}(n, \lambda/n)$, then

$$\begin{aligned} \mathbb{P}(X_n = x) &= \binom{n}{x} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \\ &= \frac{\lambda^x}{x!} \frac{n \times (n-1) \times \dots \times (n-x+1)}{n \times n \times \dots \times n} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}. \end{aligned}$$

As $n \rightarrow \infty$ the second and fourth factors converge to 1 and the third factor to $e^{-\lambda}$ (this is one of the defining properties of the exponential function). Hence, we have

$$\lim_{n \rightarrow \infty} \mathbb{P}(X_n = x) = \frac{\lambda^x}{x!} e^{-\lambda}.$$

An example of the graph of the Poisson pdf is given in Fig. 2.8. Theorem 2.12 summarizes some properties of the Poisson distribution.

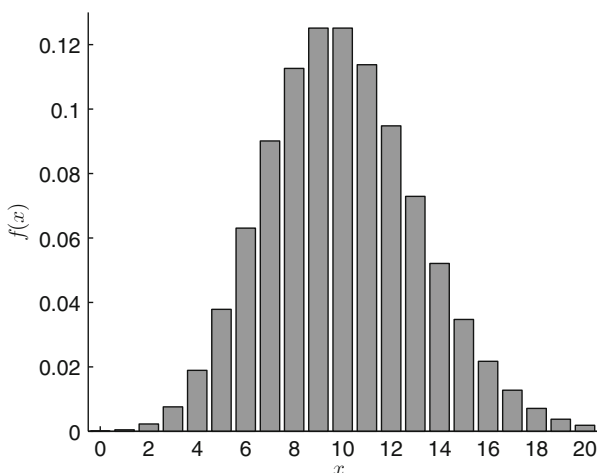



Fig. 2.8 The pdf of the $\text{Poi}(10)$ distribution

Theorem 2.12. (Properties of the Poisson Distribution). Let $X \sim \text{Poi}(\lambda)$. Then,

1. $\mathbb{E}X = \lambda$,
2. $\text{Var}(X) = \lambda$,
3. the PGF is $G(z) = e^{-\lambda(1-z)}$.

 34 *Proof.* The PGF was derived in Example 2.6. It follows from Corollary 2.2 that $\mathbb{E}X = G'(1) = \lambda$ and

$$\text{Var}(X) = G''(1) + G'(1) - (G'(1))^2 = \lambda^2 + \lambda - \lambda^2 = \lambda.$$

Thus, the parameter λ can be interpreted as both the expectation and variance of X . □

2.6 Common Continuous Distributions

In this section we give a number of common continuous distributions and list some of their properties. Note that the pdf of each of these distributions depends on one or more parameters; so, as in the previous section, we are dealing with *families* of distributions.

2.6.1 Uniform Distribution

Definition 2.14. (Uniform Distribution). A random variable X is said to have a **uniform** distribution on the interval $[a, b]$ if its pdf is given by

$$f(x) = \frac{1}{b-a}, \quad a \leq x \leq b.$$

We write $X \sim \text{U}[a, b]$ (and $X \sim \text{U}(a, b)$ for a uniform random variable on an open interval (a, b)).

Fig. 2.9 The pdf of the uniform distribution on $[a, b]$



The random variable $X \sim \mathcal{U}[a, b]$ can model a randomly chosen point from the interval $[a, b]$, where each choice is equally likely. A graph of the pdf is given in Fig. 2.9.

Theorem 2.13. (Properties of the Uniform Distribution). Let $X \sim \mathcal{U}[a, b]$. Then,

1. $\mathbb{E}X = (a + b)/2$,
2. $\text{Var}(X) = (b - a)^2/12$.

Proof. We have

$$\mathbb{E}X = \int_a^b \frac{x}{b-a} dx = \frac{1}{b-a} \left[\frac{b^2 - a^2}{2} \right] = \frac{a+b}{2}$$

and

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}X^2 - (\mathbb{E}X)^2 = \int_a^b \frac{x^2}{b-a} dx - \left(\frac{a+b}{2} \right)^2 \\ &= \frac{b^3 - a^3}{3(b-a)} - \left(\frac{a+b}{2} \right)^2 = \frac{(b-a)^2}{12}. \end{aligned}$$

□

2.6.2 Exponential Distribution

Definition 2.15. (Exponential Distribution). A random variable X is said to have an **exponential** distribution with parameter λ if its pdf is given by

$$f(x) = \lambda e^{-\lambda x}, \quad x \geq 0. \quad (2.15)$$

We write $X \sim \text{Exp}(\lambda)$.

The exponential distribution can be viewed as a continuous version of the geometric distribution. Graphs of the pdf for various values of λ are given in Fig. 2.10. Theorem 2.14 summarizes some properties of the exponential distribution.

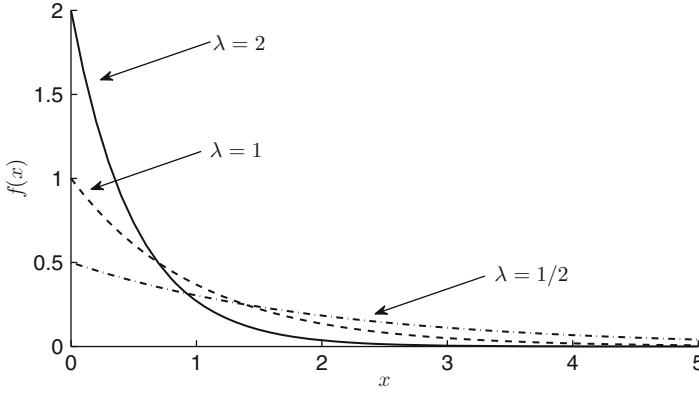


Fig. 2.10 The pdf of the $\text{Exp}(\lambda)$ distribution for various λ

Theorem 2.14. (Properties of the Exponential Distribution). Let $X \sim \text{Exp}(\lambda)$. Then,

1. $\mathbb{E}X = 1/\lambda$,
2. $\text{Var}(X) = 1/\lambda^2$,
3. the MGF of X is $M(s) = \lambda/(\lambda - s)$, $s < \lambda$,
4. the cdf of X is $F(x) = 1 - e^{-\lambda x}$, $x \geq 0$,
5. the **memoryless property** holds: for any $s, t > 0$,

$$\mathbb{P}(X > s + t \mid X > s) = \mathbb{P}(X > t). \quad (2.16)$$

Proof.

3. The MGF is given by

$$\begin{aligned} M(s) &= \int_0^\infty e^{sx} \lambda e^{-\lambda x} dx = \lambda \int_0^\infty e^{-(\lambda-s)x} dx = \lambda \left[\frac{-e^{-(\lambda-s)x}}{\lambda-s} \right]_0^\infty \\ &= \frac{\lambda}{\lambda-s}, \quad s < \lambda \quad (\text{and } M(s) = \infty \text{ for } s \geq \lambda). \end{aligned}$$

36 1. From Theorem 2.7, we obtain

$$\mathbb{E}X = M'(0) = \frac{\lambda}{(\lambda-s)^2} \Big|_{s=0} = \frac{1}{\lambda}.$$

2. Similarly, the second moment is $\mathbb{E}X^2 = M''(0) = \frac{2\lambda}{(\lambda-s)^3} \Big|_{s=0} = 2/\lambda^2$, so that the variance is

$$\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2} .$$

4. The cdf of X is given by

$$F(x) = \mathbb{P}(X \leq x) = \int_0^x \lambda e^{-\lambda u} du = [-e^{-\lambda u}]_0^x = 1 - e^{-\lambda x}, \quad x \geq 0 .$$

Note that the tail probability $\mathbb{P}(X > x)$ is exponentially decaying:

$$\mathbb{P}(X > x) = e^{-\lambda x}, \quad x \geq 0 .$$

5. Similar to the proof of the memoryless property for the geometric distribution (Theorem 2.11), we have

$$\begin{aligned} \mathbb{P}(X > s + t \mid X > s) &= \frac{\mathbb{P}(X > s + t, X > s)}{\mathbb{P}(X > s)} = \frac{\mathbb{P}(X > s + t)}{\mathbb{P}(X > s)} \\ &= \frac{e^{-\lambda(t+s)}}{e^{-\lambda s}} = e^{-\lambda t} = \mathbb{P}(X > t) . \end{aligned}$$

□

The memoryless property can be interpreted as a “non-aging” property. For example, when X denotes the lifetime of a machine, then, given the fact that the machine is alive at time s , the remaining lifetime of the machine, $X - s$, has the same exponential distribution as a completely new machine. In other words, the machine has no memory of its age and does not deteriorate (although it will break down eventually).

2.6.3 Normal (Gaussian) Distribution

In this section we introduce the most important distribution in the study of statistics: the normal (or Gaussian) distribution. Additional properties of this distribution will be given in Sect. 3.6.

Definition 2.16. (Normal Distribution). A random variable X is said to have a **normal** distribution with parameters μ and σ^2 if its pdf is given by

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, \quad x \in \mathbb{R} . \quad (2.17)$$

We write $X \sim \mathcal{N}(\mu, \sigma^2)$.

The parameters μ and σ^2 turn out to be the expectation and variance of the distribution, respectively. If $\mu = 0$ and $\sigma = 1$, then

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2},$$

and the distribution is known as the **standard normal** distribution. The cdf of the standard normal distribution is often denoted by Φ and its pdf by φ . In Fig. 2.11 the pdf of the $N(\mu, \sigma^2)$ distribution is plotted for various μ and σ^2 .

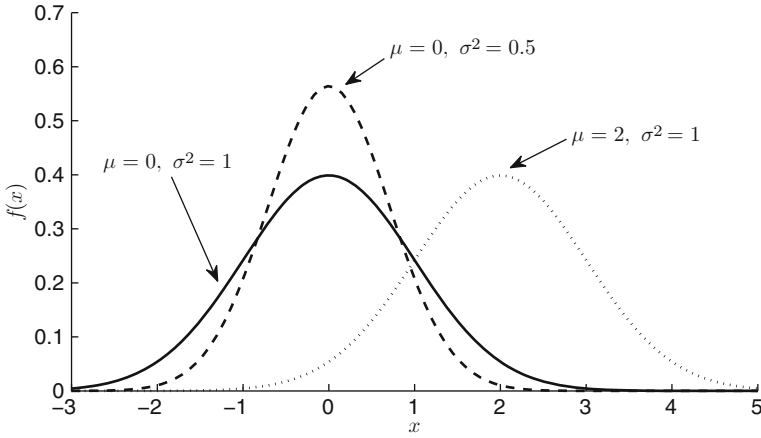


Fig. 2.11 The pdf of the $N(\mu, \sigma^2)$ distribution for various μ and σ^2

We next consider some important properties of the normal distribution.

Theorem 2.15. (Standardization). Let $X \sim N(\mu, \sigma^2)$ and define $Z = (X - \mu)/\sigma$. Then Z has a standard normal distribution.

Proof. The cdf of Z is given by

$$\begin{aligned} \mathbb{P}(Z \leq z) &= \mathbb{P}((X - \mu)/\sigma \leq z) = \mathbb{P}(X \leq \mu + \sigma z) \\ &= \int_{-\infty}^{\mu + \sigma z} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy = \Phi(z), \end{aligned}$$

where we make a change of variable $y = (x - \mu)/\sigma$ in the fourth equation. Hence, $Z \sim N(0, 1)$. \square

The rescaling procedure in Theorem 2.15 is called **standardization**. It follows from Theorem 2.15 that any $X \sim N(\mu, \sigma^2)$ can be written as

$$X = \mu + \sigma Z, \quad \text{where } Z \sim N(0, 1).$$

In other words, any normal random variable can be viewed as an **affine transformation**—that is, a linear transformation plus a constant—of a standard normal random variable.

Next we prove the earlier claim that the parameters μ and σ^2 are, respectively, the expectation and variance of the distribution.

Theorem 2.16. (Expectation and Variance for the Normal Distribution).

If $X \sim N(\mu, \sigma^2)$, then $\mathbb{E}X = \mu$ and $\text{Var}(X) = \sigma^2$.

Proof. Since the pdf is symmetric around μ and $\mathbb{E}X < \infty$, it follows that $\mathbb{E}X = \mu$. To show that the variance of X is σ^2 , we first write $X = \mu + \sigma Z$, where $Z \sim N(0, 1)$. Then, $\text{Var}(X) = \text{Var}(\mu + \sigma Z) = \sigma^2 \text{Var}(Z)$. Hence, it suffices to show that $\text{Var}(Z) = 1$. Now, since the expectation of Z is 0, we have

$$\text{Var}(Z) = \mathbb{E}Z^2 = \int_{-\infty}^{\infty} z^2 \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = \int_{-\infty}^{\infty} z \times \frac{z}{\sqrt{2\pi}} e^{-z^2/2} dz.$$

We apply integration by parts to the last integral to find

$$\mathbb{E}Z^2 = \left[-\frac{z}{\sqrt{2\pi}} e^{-z^2/2} \right]_{-\infty}^{\infty} + \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = 1,$$

since the last integrand is the pdf of the standard normal distribution. \square

Theorem 2.17. (MGF for the Normal Distribution). The MGF of $X \sim N(\mu, \sigma^2)$ is

$$\mathbb{E}e^{sX} = e^{s\mu + s^2\sigma^2/2}, \quad s \in \mathbb{R}. \quad (2.18)$$

Proof. Write $X = \mu + \sigma Z$, where $Z \sim N(0, 1)$. We have

$$\mathbb{E}e^{sZ} = \int_{-\infty}^{\infty} e^{sz} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = e^{s^2/2} \underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-(z-s)^2/2} dz}_{\text{pdf of } N(s, 1)} = e^{s^2/2},$$

so that $\mathbb{E}e^{sX} = \mathbb{E}e^{s(\mu + \sigma Z)} = e^{s\mu} \mathbb{E}e^{s\sigma Z} = e^{s\mu} e^{\sigma^2 s^2/2} = e^{s\mu + \sigma^2 s^2/2}$. \square

2.6.4 Gamma and χ^2 Distribution

Definition 2.17. (Gamma Distribution). A random variable X is said to have a **gamma** distribution with **shape** parameter $\alpha > 0$ and **scale** parameter $\lambda > 0$ if its pdf is given by

$$f(x) = \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)}, \quad x \geq 0, \quad (2.19)$$

where Γ is the gamma function. We write $X \sim \text{Gamma}(\alpha, \lambda)$.

The **gamma function** $\Gamma(\alpha)$ is an important special function in mathematics, defined by

$$\Gamma(\alpha) = \int_0^\infty u^{\alpha-1} e^{-u} du. \quad (2.20)$$

We mention a few properties of the Γ function:

1. $\Gamma(\alpha + 1) = \alpha \Gamma(\alpha)$ for $\alpha \in \mathbb{R}_+$.
2. $\Gamma(n) = (n - 1)!$ for $n = 1, 2, \dots$
3. $\Gamma(1/2) = \sqrt{\pi}$.

Two special cases of the $\text{Gamma}(\alpha, \lambda)$ distribution are worth mentioning. Firstly, the $\text{Gamma}(1, \lambda)$ distribution is simply the $\text{Exp}(\lambda)$ distribution. Secondly, the $\text{Gamma}(n/2, 1/2)$ distribution, where $n \in \{1, 2, \dots\}$, is called the **chi-squared** distribution with n **degrees of freedom**. We write $X \sim \chi_n^2$. A graph of the pdf of the χ_n^2 distribution for various n is given in Fig. 2.12.

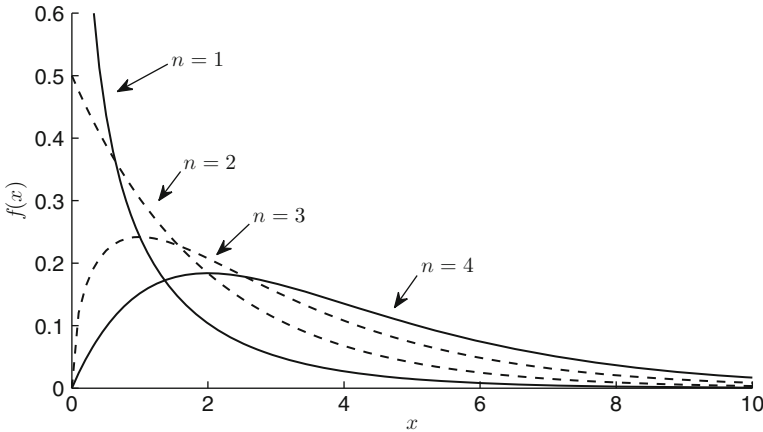


Fig. 2.12 The pdf of the χ_n^2 distribution for various degrees of freedom n

The following theorem summarizes some properties of the gamma distribution.

Theorem 2.18. (Properties of the Gamma Distribution). Let $X \sim \text{Gamma}(\alpha, \lambda)$. Then,

1. $\mathbb{E}X = \alpha/\lambda$,
2. $\text{Var}(X) = \alpha/\lambda^2$,
3. the MGF is $M(s) = [\lambda/(\lambda - s)]^\alpha$, $s < \lambda$ (and ∞ otherwise).

Proof.

3. For $s < \lambda$, the MGF of X at s is given by

$$\begin{aligned}
 M(s) = \mathbb{E} e^{sX} &= \int_0^\infty \frac{e^{-\lambda x} \lambda^\alpha x^{\alpha-1}}{\Gamma(\alpha)} e^{sx} dx \\
 &= \left(\frac{\lambda}{\lambda - s} \right)^\alpha \int_0^\infty \underbrace{\frac{e^{-(\lambda-s)x} (\lambda-s)^\alpha x^{\alpha-1}}{\Gamma(\alpha)}}_{\text{pdf of Gamma}(\alpha, \lambda-s)} dx \\
 &= \left(\frac{\lambda}{\lambda - s} \right)^\alpha.
 \end{aligned} \tag{2.21}$$

1. Consequently, by Theorem 2.7,

 36

$$\mathbb{E}X = M'(0) = \frac{\alpha}{\lambda} \left(\frac{\lambda}{\lambda - s} \right)^{\alpha+1} \Big|_{s=0} = \frac{\alpha}{\lambda}.$$

2. Similarly, $\text{Var}(X) = M''(0) - (M'(0))^2 = (\alpha + 1)\alpha/\lambda^2 - (\alpha/\lambda)^2 = \alpha/\lambda^2$. \square

2.6.5 F Distribution

Definition 2.18. (F Distribution). Let m and n be strictly positive integers. A random variable X is said to have an **F** distribution with **degrees of freedom** m and n if its pdf is given by

$$f(x) = \frac{\Gamma(\frac{m+n}{2}) (m/n)^{m/2} x^{(m-2)/2}}{\Gamma(\frac{m}{2}) \Gamma(\frac{n}{2}) [1 + (m/n)x]^{(m+n)/2}}, \quad x \geq 0, \tag{2.22}$$

where Γ denotes the gamma function. We write $X \sim F(m, n)$.

88 The F distribution plays an important role in classical statistics, through Theorem 3.11. A graph of the pdf of the $F(m, n)$ distribution for various m and n is given in Fig. 2.13.

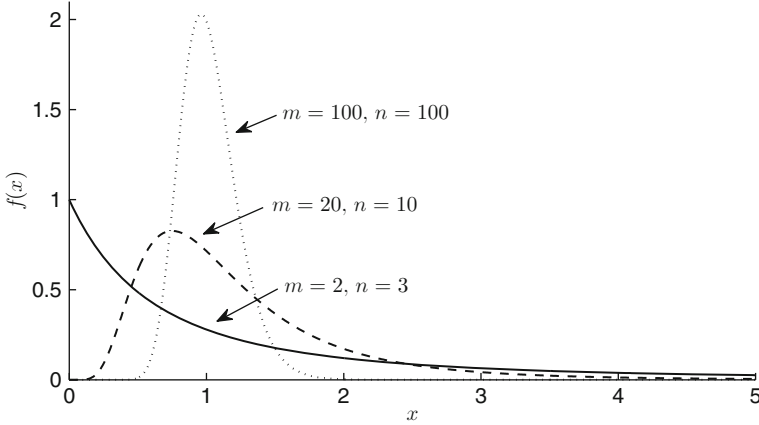


Fig. 2.13 The pdf of the $F(m, n)$ distribution for various degrees of freedom m and n

2.6.6 Student's t Distribution

Definition 2.19. (Student's t Distribution). A random variable X is said to have a **Student's t** distribution with parameter $\nu > 0$ if its pdf is given by

$$f(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-(\nu+1)/2}, \quad x \in \mathbb{R}, \quad (2.23)$$

where Γ denotes the gamma function. We write $X \sim t_\nu$. For integer values the parameter ν is referred to as the **degrees of freedom** of the distribution.

A graph of the pdf of the t_ν distribution for various ν is given in Fig. 2.14. Note that the pdf is symmetric. Moreover, it can be shown that the pdf of the t_ν distribution converges to the pdf of the $N(0, 1)$ distribution as $\nu \rightarrow \infty$. The t_1 distribution is called the **Cauchy distribution**.

For completeness we mention that if $X \sim t_\nu$, then

$$\mathbb{E}X = 0 \quad (\nu > 1) \quad \text{and} \quad \text{Var}(X) = \frac{\nu}{\nu - 2}, \quad (\nu > 2).$$

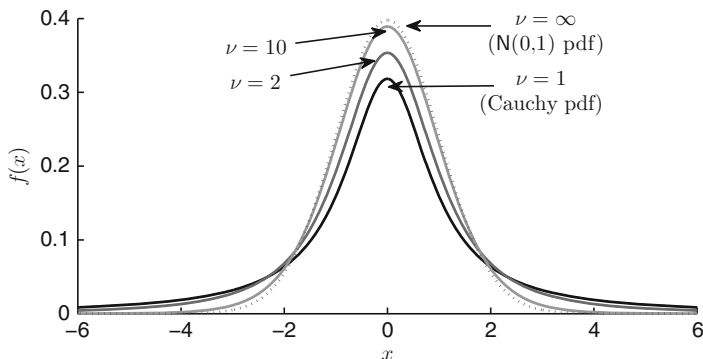


Fig. 2.14 The pdfs of t_1 (Cauchy), t_2 , t_{10} , and t_∞ ($N(0, 1)$) distributions

The t and F distributions are related in the following way.

Theorem 2.19. (Relationship Between the t and F Distribution). For integer $n \geq 1$, if $X \sim t_n$, then $X^2 \sim F(1, n)$.

Proof. Let $Z = X^2$. We can express the cdf of Z in terms of the cdf of X . Namely, for every $z > 0$, we have

$$F_Z(z) = \mathbb{P}(X^2 \leq z) = \mathbb{P}(-\sqrt{z} \leq X \leq \sqrt{z}) = F_X(\sqrt{z}) - F_X(-\sqrt{z}).$$

Differentiating with respect to z gives the following relation between the two pdfs:

$$f_Z(z) = f_X(\sqrt{z}) \frac{1}{2\sqrt{z}} + f_X(-\sqrt{z}) \frac{1}{2\sqrt{z}} = f_X(\sqrt{z}) \frac{1}{\sqrt{z}},$$

using the symmetry of the t distribution. Substituting (2.23) into the last equation yields

$$f_Z(z) = c(n) \frac{z^{-1/2}}{(1 + z/n)^{(n+1)/2}}, \quad z > 0$$

for some constant $c(n)$. The only pdf of this form is that of the $F(1, n)$ distribution. \square

2.7 Generating Random Variables

This section shows how to generate random variables on a computer. We first discuss a modern uniform random generator and then introduce two general methods

for drawing from an arbitrary one-dimensional distribution: the inverse-transform method and the acceptance–rejection method.

2.7.1 Generating Uniform Random Variables

The MATLAB `rand` function simulates the drawing of a uniform random number on the interval $(0, 1)$ by generating *pseudorandom* numbers, that is, numbers that, although not actually random (because the computer is a deterministic device), behave for all intended purposes as truly random. The following algorithm (L’Ecuyer 1999) uses simple recurrences to produce high-quality pseudorandom numbers, in the sense that the numbers pass all currently known statistical tests for randomness and uniformity.

Algorithm 2.1. (Combined Multiple-Recursive Generator).

1. Suppose N random numbers are required. Define $m_1 = 2^{32} - 209$ and $m_2 = 2^{32} - 22853$.
2. Initialize a vector $(X_{-2}, X_{-1}, X_0) = (12345, 12345, 12345)$ and a vector $(Y_{-2}, Y_{-1}, Y_0) = (12345, 12345, 12345)$.
3. For $t = 1$ to N let

$$X_t = (1403580 X_{t-2} - 810728 X_{t-3}) \bmod m_1 ,$$

$$Y_t = (527612 Y_{t-1} - 1370589 Y_{t-3}) \bmod m_2 ,$$

and output the t th random number as

$$U_t = \begin{cases} \frac{X_t - Y_t + m_1}{m_1 + 1} & \text{if } X_t \leq Y_t , \\ \frac{X_t - Y_t}{m_1 + 1} & \text{if } X_t > Y_t . \end{cases}$$

Here, $x \bmod m$ means the remainder of x when divided by m . The initialization in Step 2 determines the initial state—the so-called **seed**—of the random number stream. Restarting the stream from the same seed produces the same sequence.

Algorithm 2.1 is implemented as a core MATLAB uniform random number generator from Version 7. Currently the default generator in MATLAB is the *Mersenne twister*, which also passes (most) statistical tests and tends to be a little faster. However, it is considerably more difficult to implement. A typical usage of MATLAB’s uniform random number generator is as follows.

```

>>rng(1,'combRecursive') % use the CMRG with seed 1
>>rand(1,5) % draw 5 random numbers

ans =
    0.4957    0.2243    0.2073    0.6823    0.6799

>>rng(1234) % set the seed to 1234
>>rand(1,5)

ans =
    0.2830    0.2493    0.3600    0.9499    0.8071

>>rng(1234) % reset the seed to 1234

>>rand(1,5)

ans =
    0.2830    0.2493    0.3600    0.9499    0.8071

```

2.7.2 Inverse-Transform Method

Once we have a method for drawing a uniform random number, we can, in principle, simulate a random variable X from *any* cdf F by using the following algorithm.

Algorithm 2.2. (Inverse-Transform Method).

1. Generate U from $U(0, 1)$.
2. Return $X = F^{-1}(U)$, where F^{-1} is the inverse function of F .

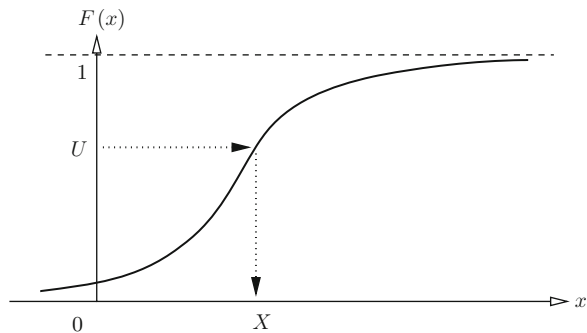


Fig. 2.15 The inverse-transform method

Figure 2.15 illustrates the inverse-transform method. We see that the random variable $X = F^{-1}(U)$ has cdf F , since

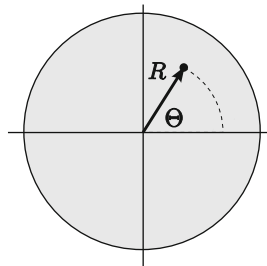
$$\mathbb{P}(X \leq x) = \mathbb{P}(F^{-1}(U) \leq x) = \mathbb{P}(U \leq F(x)) = F(x). \quad (2.24)$$

Example 2.7 (Generating Uniformly on a Unit Disk). Suppose we wish to draw a random point (X, Y) uniformly on the unit disk; see Fig. 2.16. In polar coordinates we have $X = R \cos \Theta$ and $Y = R \sin \Theta$, where Θ has a $U(0, 2\pi)$ distribution. The cdf of R is given by

$$F(r) = \mathbb{P}(R \leq r) = \frac{\pi r^2}{\pi} = r^2, \quad 0 < r < 1.$$

Its inverse is $F^{-1}(u) = \sqrt{u}$, $0 < u < 1$. We can thus generate R via the inverse-transform method as $R = \sqrt{U_1}$, where $U_1 \sim U(0, 1)$. In addition, we can simulate Θ as $\Theta = 2\pi U_2$, where $U_2 \sim U(0, 1)$. Note that U_1 and U_2 should be independent draws from $U(0, 1)$.

Fig. 2.16 Draw a point (X, Y) uniformly on the unit disk



The inverse-transform method holds for general cdfs F . Note that F for discrete random variables is a step function, as illustrated in Fig. 2.17. The algorithm for generating a random variable X from a discrete distribution that takes values x_1, x_2, \dots with probabilities p_1, p_2, \dots is thus as follows.

Algorithm 2.3. (Discrete Inverse-Transform Method).

1. Generate $U \sim U(0, 1)$.
2. Find the smallest positive integer k such that $F(x_k) \geq U$ and return $X = x_k$.

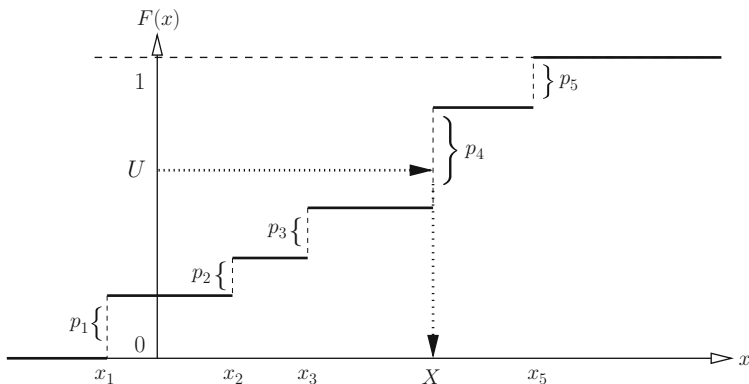


Fig. 2.17 The inverse-transform method for a discrete random variable

Drawing one of the numbers $1, \dots, n$ according to a probability vector (p_1, \dots, p_n) can be done in one line of MATLAB code:

```
min(find(cumsum(p) > rand));
```

Here \mathbf{p} is the vector of probabilities, such as $(0.3, 0.2, 0.5)$, `cumsum` gives the cumulative vector, e.g., $(0.3, 0.5, 1)$, `find(...)` finds the indices i such that the cumulative probability is greater than some random number `rand`, and `min` takes the smallest of these indices.

2.7.3 Acceptance–Rejection Method

The inverse-transform method may not always be easy to implement, in particular when the inverse cdf is difficult to compute. In that case the **acceptance–rejection** method may prove to be useful. The idea of this method is depicted in Fig. 2.18. Suppose we wish to sample from a pdf f . Let g be another pdf such that for some constant $C \geq 1$ we have that $Cg(x) \geq f(x)$ for all x . It is assumed that it is easy to sample from g , for example, via the inverse-transform method.

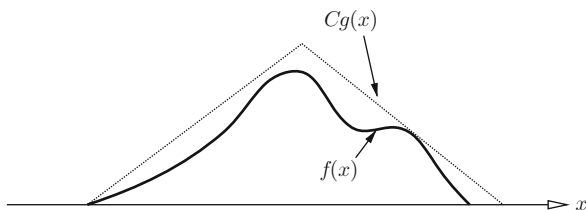


Fig. 2.18 Illustration of the acceptance–rejection method

It is intuitively clear that if a random point (X, Y) is *uniformly* distributed under the graph of f —that is, on the set $\{(x, y) : 0 \leq y \leq f(x)\}$ —then X must have pdf f . To construct such a point, let us first draw a random point (Z, V) by drawing Z from g and then drawing V uniformly on $[0, Cg(Z)]$. The point (Z, V) is uniformly distributed under the graph of Cg . If we keep drawing such a point (Z, V) *until it lies under the graph of f* , then the resulting point (X, Y) must be uniformly distributed under the graph of f and hence the X coordinate must have pdf f . This leads to the following algorithm.

Algorithm 2.4. (Acceptance–Rejection Method).

1. Generate $Z \sim g$.
2. Generate $Y \sim U(0, Cg(Z))$.
3. If $Y \leq f(Z)$, return $X = Z$; otherwise, repeat from Step 1.

Example 2.8 (Generating from the Standard Normal Distribution). To sample from the standard normal pdf via the inverse-transform method requires knowledge of the inverse of the corresponding cdf, which involves numerical integration. Instead, we can use acceptance–rejection. First, observe that the standard normal pdf is symmetric around 0. Hence, if we can generate a random variable X from the **positive normal** pdf (see Fig. 2.19),

$$f(x) = \sqrt{\frac{2}{\pi}} e^{-x^2/2}, \quad x \geq 0, \quad (2.25)$$

then we can generate a standard normal random variable by multiplying X with 1 or -1 , each with probability $1/2$. We can bound $f(x)$ by $Cg(x)$, where $g(x) = e^{-x}$ is the pdf of the $\text{Exp}(1)$ distribution. The smallest constant C such that $f(x) \leq Cg(x)$ is $\sqrt{2e/\pi}$.

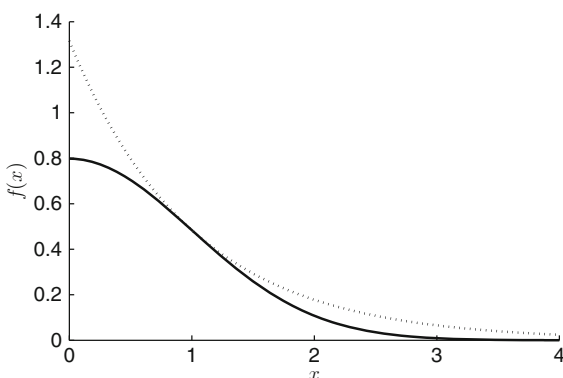


Fig. 2.19 Bounding the positive normal density (*solid curve*) via an $\text{Exp}(1)$ pdf (times $C \approx 1.3155$)

Drawing from the $\text{Exp}(1)$ distribution can be easily done via the inverse-transform method, noting that the corresponding cdf is the function $1 - e^{-x}$, $x \geq 0$, whose inverse is the function $-\ln(1 - u)$, $u \in (0, 1)$. This gives the following specification of Algorithm 2.4, where f and C are defined above.

Algorithm 2.5. (N(0, 1) Generator).

1. Draw $U_1 \sim \text{U}(0, 1)$, and let $Z = -\ln U_1$.
2. Draw $U_2 \sim \text{U}(0, 1)$, and let $Y = U_2 C e^{-Z}$.
3. If $Y \leq f(Z)$, let $X = Z$ and continue with Step 4. Otherwise, repeat from Step 1.
4. Draw $U_3 \sim \text{U}(0, 1)$ and return $\tilde{X} = X(2I_{\{U_3 < 1/2\}} - 1)$ as a standard normal random variable.

In Step 1, we have used the fact that if $U \sim \text{U}(0, 1)$, then also $1 - U \sim \text{U}(0, 1)$. In Step 4, $I_{\{U_3 < 1/2\}}$ denotes the **indicator** of the event $\{U_3 < 1/2\}$, which is 1 if $U_3 < 1/2$ and 0 otherwise. An alternative generation method is given in Algorithm 3.2. In MATLAB normal random variable generation is implemented via the `randn` function.

 82

2.8 Problems

2.1. Two fair dice are thrown and the smallest of the face values, M say, is noted.

- (a) Give the discrete pdf of M in table form, as in Table 2.1.
- (b) What is the probability that M is at least 3?
- (c) Calculate the expectation and variance of M .

 27

2.2. A continuous random variable X has cdf

$$F(x) = \begin{cases} 0, & x < 0 \\ x^2/5, & 0 \leq x \leq 1 \\ \frac{1}{5}(-x^2 + 6x - 4), & 1 < x \leq 3 \\ 1, & x > 3. \end{cases}$$

- (a) Find the corresponding pdf and plot its graph.
- (b) Calculate the following probabilities:
 - (i) $\mathbb{P}(X \leq 2)$.
 - (ii) $\mathbb{P}(1 < X \leq 2)$.
 - (iii) $\mathbb{P}(1 \leq X \leq 2)$.
 - (iv) $\mathbb{P}(X > 1/2)$.

(c) Show that $\mathbb{E}X = 22/15$.

2.3. In this book most random variables are either discrete or continuous; that is, they have either a discrete or a continuous pdf. It is also possible to define random variables that have a mix of discrete and continuous characteristics. A simple example is a random variable X with cdf

$$F(x) = \begin{cases} 0, & x < 0 \\ 1 - c e^{-x}, & x \geq 0 \end{cases}$$

for some fixed $0 < c < 1$.

- (a) Sketch the cdf F .
- (b) Find the following probabilities:
 - (i) $\mathbb{P}(0 \leq X \leq x), x \geq 0$.
 - (ii) $\mathbb{P}(0 < X \leq x), x \geq 0$.
 - (iii) $\mathbb{P}(X = x), x \geq 0$.
- (c) Describe how the inverse-transform method can be used to draw samples from this distribution.

2.4. Let X be a positive random variable with cdf F . Prove that

$$\mathbb{E}X = \int_0^\infty (1 - F(x)) dx. \quad (2.26)$$

2.5. Let X be a random variable that can possibly take values $-\infty$ and ∞ with probabilities $\mathbb{P}(X = -\infty) = a$ and $\mathbb{P}(X = \infty) = b$, respectively. Show that the corresponding cdf F satisfies $\lim_{x \rightarrow -\infty} F(x) = a$ and $\lim_{x \rightarrow \infty} F(x) = 1 - b$.

2.6. Suppose that in a large population the fraction of left-handers is 12%. We select at random 100 people from this population. Let X be the number of left-handers among the selected people. What is the distribution of X ? What is the probability that at most 7 of the selected people are left-handed?

2.7. Let $X \sim \text{Geom}(p)$. Show that

$$\mathbb{P}(X > k) = (1 - p)^k.$$

2.8. Find the MGF of $X \sim U[a, b]$.

2.9. Let $X = a + (b - a)U$, where $U \sim U[0, 1]$. Prove that $X \sim U[a, b]$. Use this to provide a more elegant proof of Theorem 2.13.

43

2.10. Show that the exponential distribution is the *only* continuous (positive) distribution that possesses the memoryless property. Hint: show that the memoryless property implies that the tail probability $g(x) = \mathbb{P}(X > x)$ satisfies $g(x + y) = g(x)g(y)$.

2.11. Let $X \sim \text{Exp}(2)$. Calculate the following quantities:

- (a) $\mathbb{P}(-1 \leq X \leq 1)$.
- (b) $\mathbb{P}(X > 4)$.
- (c) $\mathbb{P}(X > 4 \mid X > 2)$.
- (d) $\mathbb{E}X^2$.

2.12. What is the expectation of a random variable X with the following discrete pdf on the set of integer numbers, excluding 0:

$$f(x) = \frac{3}{\pi^2} \frac{1}{x^2}, \quad x \in \mathbb{Z} \setminus \{0\}.$$

What is the pdf of the absolute value $|X|$ and what is its expectation?

2.13. A random variable X is said to have a **discrete uniform distribution** on the set $\{a, a+1, \dots, b\}$ if

$$\mathbb{P}(X = x) = \frac{1}{b - a + 1}, \quad x = a, a+1, \dots, b.$$

- (a) What is the expectation of X ?
- (b) Show that $\text{Var}(X) = (b-a)(b-a+1)/12$.
- (c) Find the PGF of X .
- (d) Describe a simple way to generate X using a uniform number generator.

2.14. Let X and Y be random variables. Prove that if $X \leq Y$, then $\mathbb{E}X \leq \mathbb{E}Y$.

2.15. A continuous random variable is said to have a **logistic** distribution if its pdf is given by

$$f(x) = \frac{e^{-x}}{(1 + e^{-x})^2}, \quad x \in \mathbb{R}. \quad (2.27)$$

- (a) Plot the graph of the pdf.
- (b) Show that $\mathbb{P}(X > x) = 1/(1 + e^x)$ for all x .
- (c) Write an algorithm based on the inverse-transform method to generate random variables from this distribution.

2.16. An electrical component has a lifetime (in years) that is distributed according to an exponential distribution with expectation 3. What is the probability that the component is still functioning after 4.5 years, given that it still works after 4 years? Answer the same question for the case where the component's lifetime is normally distributed with the same expected value and variance as before.

2.17. Consider the pdf given by

$$f(x) = \begin{cases} 4e^{-4(x-1)}, & x \geq 1, \\ 0, & x < 1. \end{cases}$$

- (a) If X is distributed according to this pdf f , what is its expectation?
- (b) Specify how one can generate a random variable $X \sim f$ using a uniform random number generator.

2.18. Let $X \sim N(4, 9)$.

- (a) Plot the graph of the pdf.
- (b) Express the following probabilities in terms of the cdf Φ of the standard normal distribution:
 - (i) $\mathbb{P}(X \leq 3)$.
 - (ii) $\mathbb{P}(X > 4)$.
 - (iii) $\mathbb{P}(-1 \leq X \leq 5)$.
- (c) Find $\mathbb{E}[2X + 1]$.
- (d) Calculate $\mathbb{E}X^2$.



2.19. Let Φ be the cdf of $X \sim N(0, 1)$. The integral

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2} du$$

needs to be evaluated numerically. In MATLAB there are several ways to do this:

- (1) If the *Statistics Toolbox* is available, the cdf can be evaluated via the functions `normcdf` or `cdf`. The inverse cdf can be evaluated using `norminv` or `icdf`. See also their replacements `cumcdf` and `icumcdf` in Appendix A.9.
- (2) Or one could use the built-in **error function** `erf`, defined as

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-u^2} du, \quad x \in \mathbb{R}.$$

The inverse of the error function, erf^{-1} , is implemented in MATLAB as `erfinv`.

- (3) A third alternative is to use numerical integration (quadrature) via the `quad` function. For example, `quad(@f, 0, 1)` integrates a MATLAB function `f.m` on the interval $[0, 1]$.
- (a) Show that $\Phi(x) = (\text{erf}(x/\sqrt{2}) + 1)/2$.
- (b) Evaluate $\Phi(x)$ for $x = 1, 2$, and 3 via (a) the error function and (b) numerical integration of the pdf, using the fact that $\Phi(0) = 1/2$.
- (c) Show that the inverse of Φ is given by

$$\Phi^{-1}(y) = \sqrt{2} \text{erf}^{-1}(2y - 1), \quad 0 < y < 1.$$



2.20. Based on MATLAB's `rand` and `randn` functions *only*, implement algorithms that generate random variables from the following distributions:

- (a) $U[2, 3]$.
- (b) $N(3, 9)$.
- (c) $\text{Exp}(4)$.
- (d) $\text{Bin}(10, 1/2)$.
- (e) $\text{Geom}(1/6)$.



2.21. The **Weibull** distribution $\text{Weib}(\alpha, \lambda)$ has cdf

$$F(x) = 1 - e^{-(\lambda x)^\alpha}, \quad x \geq 0. \quad (2.28)$$

It can be viewed as a generalization of the exponential distribution. Write a **MATLAB** program that draws 1000 samples from the $\text{Weib}(2, 1)$ distribution using the inverse-transform method. Give a histogram of the sample.



2.22. Consider the pdf

$$f(x) = c e^{-x} x(1 - x), \quad 0 \leq x \leq 1.$$

- (a) Show that $c = e/(3 - e)$.
- (b) Devise an acceptance–rejection algorithm to generate random variables that are distributed according to f .
- (c) Implement the algorithm in **MATLAB**.



2.23. Implement two different algorithms to draw 100 uniformly generated points on the unit disk: one based on Example 2.7 and the other using (two-dimensional) acceptance–rejection.

Chapter 3

Joint Distributions

Often a random experiment is described via more than one random variable. Here are some examples:

1. We randomly select $n = 10$ people and observe their heights. Let X_1, \dots, X_n be the individual heights.
2. We toss a coin repeatedly. Let $X_i = 1$ if the i th toss is Heads and $X_i = 0$ otherwise. The experiment is thus described by the sequence X_1, X_2, \dots of Bernoulli random variables.
3. We randomly select a person from a large population and measure his/her weight X and height Y .

How can we specify the behavior of the random variables above? We should not just specify the pdf of the individual random variables, but also say something about the interaction (or lack thereof) between the random variables. For example, in the third experiment above, if the height Y is large, then most likely X is large as well. In contrast, in the first two experiments, it is reasonable to assume that the random variables are “independent” in some way; that is, information about one of the random variables does not give extra information about the others. What we need to specify is the **joint distribution** of the random variables. The theory below for multiple random variables follows a similar path to that of a single random variable described in Sects. 2.1–2.3.

Let X_1, \dots, X_n be random variables describing some random experiment. We can accumulate the $\{X_i\}$ into a **random vector** $\mathbf{X} = (X_1, \dots, X_n)$ (row vector) or $\mathbf{X} = (X_1, \dots, X_n)^\top$ (column vector). Recall that the distribution of a *single* random variable X is completely specified by its cumulative distribution function. For *multiple* random variables we have the following generalization.

Definition 3.1. (Joint Cdf). The **joint cdf** of X_1, \dots, X_n is the function $F : \mathbb{R}^n \rightarrow [0, 1]$ defined by

$$F(x_1, \dots, x_n) = \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n) .$$

 23

Notice that we have used the abbreviation $\mathbb{P}(\{X_1 \leq x_1\} \cap \cdots \cap \{X_n \leq x_n\}) = \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n)$ to denote the probability of the intersection of events. We will use this abbreviation throughout the book.

As in the univariate (i.e., single variable) case we distinguish between *discrete* and *continuous* distributions.

3.1 Discrete Joint Distributions

Example 3.1 (Dice Experiment). In a box there are three dice. Die 1 is an ordinary die; die 2 has no 6 face, but instead two 5 faces; die 3 has no 5 face, but instead two 6 faces. The experiment consists of selecting a die at random followed by a toss with that die. Let X be the die number that is selected and let Y be the face value of that die. The probabilities $\mathbb{P}(X = x, Y = y)$ in Table 3.1 specify the joint distribution of X and Y . Note that it is more convenient to specify the joint probabilities $\mathbb{P}(X = x, Y = y)$ than the joint cumulative probabilities $\mathbb{P}(X \leq x, Y \leq y)$. The latter can be found, however, from the former by applying the sum rule. For example, $\mathbb{P}(X \leq 2, Y \leq 3) = \mathbb{P}(X = 1, Y = 1) + \cdots + \mathbb{P}(X = 2, Y = 3) = 6/18 = 1/3$. Moreover, by that same sum rule, the distribution of X is found by summing the $\mathbb{P}(X = x, Y = y)$ over all values of y —giving the last column of Table 3.1. Similarly, the distribution of Y is given by the column totals in the last row of the table.

Table 3.1 The joint distribution of X (die number) and Y (face value)

		y						
		1	2	3	4	5	6	Σ
x	1	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{3}$
	2	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{9}$	0	$\frac{1}{3}$
	3	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	$\frac{1}{18}$	0	$\frac{1}{9}$	$\frac{1}{3}$
Σ		$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	1

In general, for discrete random variables X_1, \dots, X_n , the joint distribution is easiest to specify via the joint pdf.

Definition 3.2. (Discrete Joint Pdf). The **joint pdf** f of discrete random variables X_1, \dots, X_n is given by the function

$$f(x_1, \dots, x_n) = \mathbb{P}(X_1 = x_1, \dots, X_n = x_n) .$$

We sometimes write f_{X_1, \dots, X_n} instead of f to show that this is the pdf of the random variables X_1, \dots, X_n . Or, if $\mathbf{X} = (X_1, \dots, X_n)$ is the corresponding random vector, we can write $f_{\mathbf{X}}$ instead.

If the joint pdf f is known, we can calculate the probability of any event $\{\mathbf{X} \in B\}$, $B \subset \mathbb{R}^n$, via the sum rule as

$$\mathbb{P}(\mathbf{X} \in B) = \sum_{\mathbf{x} \in B} f(\mathbf{x}) .$$

Compare this with (2.2). In particular, as explained in Example 3.1, we can find the pdf of X_i —often referred to as a **marginal** pdf, to distinguish it from the joint pdf—by summing the joint pdf over all possible values of the other variables:

$$\mathbb{P}(X_i = x) = \sum_{x_1} \cdots \sum_{x_{i-1}} \sum_{x_{i+1}} \cdots \sum_{x_n} f(x_1, \dots, x_{i-1}, x, x_{i+1}, x_n) . \quad (3.1)$$

The converse is not true: from the marginal distributions one cannot in general reconstruct the joint distribution. For example, in Example 3.1, we cannot reconstruct the inside of the two-dimensional table if only given the column and row totals.

However, there is an important exception, namely, when the random variables are *independent*. We have so far only defined what independence is for *events*. We can define random variables X_1, \dots, X_n to be independent if events $\{X_1 \in B_1\}, \dots, \{X_n \in B_n\}$ are independent for any choice of sets $\{B_i\}$. Intuitively, this means that any information about one of the random variables does not affect our knowledge about the others.

Definition 3.3. (Independence). Random variables X_1, \dots, X_n are called **independent** if for all events $\{X_i \in B_i\}$ with $B_i \subset \mathbb{R}$, $i = 1, \dots, n$

$$\mathbb{P}(X_1 \in B_1, \dots, X_n \in B_n) = \mathbb{P}(X_1 \in B_1) \cdots \mathbb{P}(X_n \in B_n) . \quad (3.2)$$

A direct consequence of the above definition is the following important theorem.

Theorem 3.1. (Independence and Product Rule). Random variables X_1, \dots, X_n with joint pdf f are independent if and only if

$$f(x_1, \dots, x_n) = f_{X_1}(x_1) \cdots f_{X_n}(x_n) \quad (3.3)$$

for all x_1, \dots, x_n , where $\{f_{X_i}\}$ are the marginal pdfs.

Proof. The theorem is true in both the discrete and continuous case. We only show the discrete case, where (3.3) is a special case of (3.2). It follows that (3.3) is a *necessary* condition for independence. To see that it is also a *sufficient* condition, let $\mathbf{X} = (X_1, \dots, X_n)$ and observe that

$$\begin{aligned} \mathbb{P}(X_1 \in B_1, \dots, X_n \in B_n) &= \mathbb{P}(\mathbf{X} \in \underbrace{B_1 \times \dots \times B_n}_A) = \sum_{\mathbf{x} \in A} f(\mathbf{x}) \\ &= \sum_{\mathbf{x} \in A} f_{X_1}(x_1) \cdots f_{X_n}(x_n) = \sum_{x_1 \in B_1} f_{X_1}(x_1) \cdots \sum_{x_n \in B_n} f_{X_n}(x_n) \\ &= \mathbb{P}(X_1 \in B_1) \cdots \mathbb{P}(X_n \in B_n). \end{aligned}$$

Here $A = B_1 \times \dots \times B_n$ denotes the Cartesian product of B_1, \dots, B_n . \square

Example 3.2 (Dice Experiment Continued). We repeat the experiment in Example 3.1 with three ordinary fair dice. Since the events $\{X = x\}$ and $\{Y = y\}$ are now independent, each entry in the pdf table is $\frac{1}{3} \times \frac{1}{6}$. Clearly in the first experiment not *all* events $\{X = x\}$ and $\{Y = y\}$ are independent.

Remark 3.1. An *infinite* sequence X_1, X_2, \dots of random variables is said to be *independent* if for any finite choice of positive integers i_1, i_2, \dots, i_n (none of them the same) the random variables X_{i_1}, \dots, X_{i_n} are independent. Many statistical models involve random variables X_1, X_2, \dots that are **independent and identically distributed**, abbreviated as **iid**. We will use this abbreviation throughout this book and write the corresponding model as

$$X_1, X_2, \dots \stackrel{\text{iid}}{\sim} \text{Dist (or } f \text{ or } F),$$

where Dist is the common distribution with pdf f and cdf F .

Example 3.3 (Bernoulli Process). Consider the experiment where we toss a biased coin n times, with probability p of Heads. We can model this experiment in the following way. For $i = 1, \dots, n$ let X_i be the result of the i th toss: $\{X_i = 1\}$ means Heads (or success), and $\{X_i = 0\}$ means Tails (or failure). Also, let

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

Finally, assume that X_1, \dots, X_n are *independent*. The sequence

$$X_1, X_2, \dots \stackrel{\text{iid}}{\sim} \text{Ber}(p)$$

is called a **Bernoulli process** with success probability p . Let $X = X_1 + \dots + X_n$ be the total number of successes in n trials (tosses of the coin). Denote by B_k the set of all binary vectors $\mathbf{x} = (x_1, \dots, x_n)$ such that $\sum_{i=1}^n x_i = k$. Note that B_k has $\binom{n}{k}$ elements. We have for every $k = 0, \dots, n$,

$$\begin{aligned}
\mathbb{P}(X = k) &= \sum_{\mathbf{x} \in B_k} \mathbb{P}(X_1 = x_1, \dots, X_n = x_n) \\
&= \sum_{\mathbf{x} \in B_k} \mathbb{P}(X_1 = x_1) \cdots \mathbb{P}(X_n = x_n) = \sum_{\mathbf{x} \in B_k} p^k (1-p)^{n-k} \\
&= \binom{n}{k} p^k (1-p)^{n-k}.
\end{aligned}$$

In other words, $X \sim \text{Bin}(n, p)$. Compare this with Example 2.2. 👉 24

For the joint pdf of *dependent* discrete random variables we can write, as a consequence of the product rule (1.5), 👉 14

$$\begin{aligned}
f(x_1, \dots, x_n) &= \mathbb{P}(X_1 = x_1, \dots, X_n = x_n) \\
&= \mathbb{P}(X_1 = x_1) \mathbb{P}(X_2 = x_2 | X_1 = x_1) \times \cdots \\
&\quad \cdots \times \mathbb{P}(X_n = x_n | X_1 = x_1, \dots, X_{n-1} = x_{n-1}),
\end{aligned}$$

assuming that all probabilities $\mathbb{P}(X = x_1), \dots, \mathbb{P}(X_1 = x_1, \dots, X_{n-1} = x_{n-1})$ are nonzero. The function which maps, *for a fixed* x_1 , each variable x_2 to the conditional probability

$$\mathbb{P}(X_2 = x_2 | X_1 = x_1) = \frac{\mathbb{P}(X_1 = x_1, X_2 = x_2)}{\mathbb{P}(X_1 = x_1)} \quad (3.4)$$

is called the **conditional pdf** of X_2 given $X_1 = x_1$. We write it as $f_{X_2 | X_1}(x_2 | x_1)$. Similarly, the function $x_n \mapsto \mathbb{P}(X_n = x_n | X_1 = x_1, \dots, X_{n-1} = x_{n-1})$ is the conditional pdf of X_n given $X_1 = x_1, \dots, X_{n-1} = x_{n-1}$, which is written as $f_{X_n | X_1, \dots, X_{n-1}}(x_n | x_1, \dots, x_{n-1})$.

Example 3.4 (Generating Uniformly on a Triangle). We uniformly select a point (X, Y) from the triangle $T = \{(x, y) : x, y \in \{1, \dots, 6\}, y \leq x\}$ in Fig. 3.1.

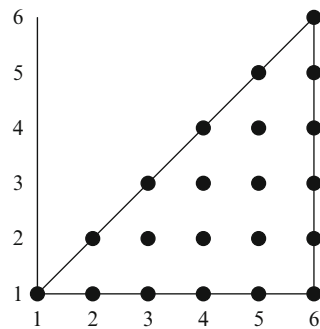


Fig. 3.1 Uniformly select a point from the triangle

Because each of the 21 points is equally likely to be selected, the joint pdf is constant on T :

$$f(x, y) = \frac{1}{21}, \quad (x, y) \in T.$$

The pdf of X is found by summing $f(x, y)$ over all y . Hence,

$$f_X(x) = \frac{x}{21}, \quad x \in \{1, \dots, 6\}.$$

Similarly,

$$f_Y(y) = \frac{7-y}{21}, \quad y \in \{1, \dots, 6\}.$$

For a fixed $x \in \{1, \dots, 6\}$ the conditional pdf of Y given $X = x$ is

$$f_{Y|X}(y|x) = \frac{f(x, y)}{f_X(x)} = \frac{1/21}{x/21} = \frac{1}{x}, \quad y \in \{1, \dots, x\},$$

which simply means that, given $X = x$, Y has a discrete uniform distribution on $\{1, \dots, x\}$.

3.1.1 Multinomial Distribution

An important discrete joint distribution is the multinomial distribution. It can be viewed as a generalization of the binomial distribution. We give the definition and then an example of how this distribution arises in applications.

Definition 3.4. (Multinomial Distribution). A random vector (X_1, X_2, \dots, X_k) is said to have a **multinomial** distribution with parameters n and p_1, p_2, \dots, p_k (positive and summing up to 1), if

$$\mathbb{P}(X_1 = x_1, \dots, X_k = x_k) = \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k} \quad (3.5)$$

for all $x_1, \dots, x_k \in \{0, 1, \dots, n\}$ such that $x_1 + x_2 + \dots + x_k = n$. We write $(X_1, \dots, X_k) \sim \text{Mnom}(n, p_1, \dots, p_k)$.

Example 3.5 (Urn Problem). We independently throw n balls into k urns, such that each ball is thrown in urn i with probability p_i , $i = 1, \dots, k$; see Fig. 3.2.

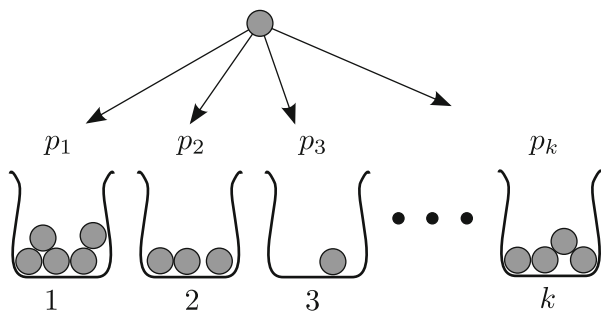


Fig. 3.2 Throwing n balls into k urns with probabilities p_1, \dots, p_k . The random configuration of balls has a multinomial distribution

Let X_i be the total number of balls in urn i , $i = 1, \dots, k$. We show that $(X_1, \dots, X_k) \sim \text{Mnom}(n, p_1, \dots, p_k)$. Let x_1, \dots, x_k be integers between 0 and n that sum up to n . The probability that the *first* x_1 balls fall in the first urn, the *next* x_2 balls fall in the second urn, etc., is

$$p_1^{x_1} p_2^{x_2} \cdots p_k^{x_k}.$$

To find the probability that there are x_1 balls in the first urn, x_2 in the second, and so on, we have to multiply the probability above with the number of ways in which we can fill the urns with x_1, x_2, \dots, x_k balls, i.e., $n!/(x_1!x_2! \cdots x_k!)$. This gives (3.5).

Remark 3.2. Note that for the *binomial* distribution there are only *two* possible urns. Also, note that for each $i = 1, \dots, k$, $X_i \sim \text{Bin}(n, p_i)$.

3.2 Continuous Joint Distributions

Joint distributions for continuous random variables are usually defined via their joint pdf. The theoretical development below follows very similar lines to both the univariate continuous case in Sect. 2.2.2 and the multivariate discrete case in Sect. 3.1.

28

64

Definition 3.5. (Continuous Joint Pdf). Continuous random variables X_1, \dots, X_n are said to have a **joint pdf** f if

$$\mathbb{P}(a_1 < X_1 \leq b_1, \dots, a_n < X_n \leq b_n) = \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} f(x_1, \dots, x_n) dx_1 \cdots dx_n$$

for all a_1, \dots, b_n .

☞ 28 This implies, similar to the univariate case in (2.3), that the probability of any event pertaining to $\mathbf{X} = (X_1, \dots, X_n)$ —say event $\{\mathbf{X} \in B\}$, where B is some subset of \mathbb{R}^n —can be found by *integration*:

$$\mathbb{P}(\mathbf{X} \in B) = \int_B f(x_1, \dots, x_n) \, dx_1 \dots dx_n . \quad (3.6)$$

☞ 29 As in (2.5) we can interpret $f(x_1, \dots, x_n)$ as the *density* of the probability distribution at (x_1, \dots, x_n) . For example, in the two-dimensional case, for small $h > 0$,

$$\begin{aligned} \mathbb{P}(x_1 \leq X_1 \leq x_1 + h, x_2 \leq X_2 \leq x_2 + h) \\ = \int_{x_1}^{x_1+h} \int_{x_2}^{x_2+h} f(u, v) \, du \, dv \approx h^2 f(x_1, x_2) . \end{aligned}$$

Similar to the discrete multivariate case in (3.1), the marginal pdfs can be recovered from the joint pdf by integrating out the other variables:

$$f_{X_i}(x) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x_1, \dots, x_{i-1}, x, x_{i+1}, \dots, x_n) \, dx_1 \dots dx_{i-1} \, dx_{i+1} \dots dx_n .$$

We illustrate this for the two-dimensional case. We have

$$F_{X_1}(x) = \mathbb{P}(X_1 \leq x, X_2 \leq \infty) = \int_{-\infty}^x \left(\int_{-\infty}^{\infty} f(x_1, x_2) \, dx_2 \right) dx_1 .$$

By differentiating the last integral with respect to x , we obtain

$$f_{X_1}(x) = \int_{-\infty}^{\infty} f(x, x_2) \, dx_2 .$$

It is not possible, in general, to reconstruct the joint pdf from the marginal pdfs. An exception is when the random variables are *independent*; see Definition 3.3. By modifying the arguments in the proof of Theorem 3.3 to the continuous case—basically replacing sums with integrals—it is not difficult to see that the theorem also holds in the continuous case. In particular, continuous random variables X_1, \dots, X_n are independent if and only if their joint pdf, f say, is the product of the marginal pdfs:

$$f(x_1, \dots, x_n) = f_{X_1}(x_1) \dots f_{X_n}(x_n) \quad (3.7)$$

for all x_1, \dots, x_n . Independence for an infinite sequence of random variables is discussed in Remark 3.1.

☞ 66

Example 3.6 (Generating a General iid Sample). Consider the sequence of numbers produced by a uniform random number generator such as MATLAB's `rand` function. A mathematical model for the output stream is U_1, U_2, \dots , are independent and $U(0, 1)$ distributed; that is,

$$U_1, U_2, \dots \stackrel{\text{iid}}{\sim} U(0, 1) .$$

Using the inverse-transform method it follows that for any cdf F ,

$$F^{-1}(U_1), F^{-1}(U_2), \dots \stackrel{\text{iid}}{\sim} F .$$

Example 3.7 (Quotient of Two Independent Random Variables). Let X and Y be independent continuous random variables, with $Y > 0$. What is the pdf of the quotient $U = X/Y$ in terms of the pdfs of X and Y ? Consider first the cdf of U . We have

$$\begin{aligned} F_U(u) &= \mathbb{P}(U \leq u) = \mathbb{P}(X/Y \leq u) = \mathbb{P}(X \leq Yu) \\ &= \int_0^\infty \int_{-\infty}^{yu} f_X(x) f_Y(y) dx dy = \int_{-\infty}^u \int_0^\infty y f_X(yz) f_Y(y) dy dz , \end{aligned}$$

where we have used the change of variable $z = x/y$ and changed the order of integration in the last equation. It follows that the pdf is given by

$$f_U(u) = \frac{d}{du} F_U(u) = \int_0^\infty y f_X(yu) f_Y(y) dy . \quad (3.8)$$

As a particular example, suppose that X and V both have a standard normal distribution. Note that X/V has the same distribution as $U = X/Y$, where $Y = |V| > 0$ has a *positive normal* distribution. It follows from (3.8) that

$$\begin{aligned} f_U(u) &= \int_0^\infty y \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2u^2} \frac{2}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} dy \\ &= \int_0^\infty y \frac{1}{\pi} e^{-\frac{1}{2}y^2(1+u^2)} dy = \frac{1}{\pi} \frac{1}{1+u^2}, \quad u \in \mathbb{R} . \end{aligned}$$

This is the pdf of the *Cauchy* distribution.

Definition 3.6. (Conditional Pdf). Let X and Y have joint pdf f and suppose $f_X(x) > 0$. The **conditional pdf** of Y given $X = x$ is defined as

$$f_{Y|X}(y|x) = \frac{f(x, y)}{f_X(x)} \quad \text{for all } y . \quad (3.9)$$

For the discrete case, this is just a rewrite of (3.4). For the continuous case, the interpretation is that $f_{Y|X}(y|x)$ is the density corresponding to the cdf $F_{Y|X}(y|x)$ defined by the limit

$$F_{Y|X}(y|x) = \lim_{h \downarrow 0} \mathbb{P}(Y \leq y | x \leq X \leq x+h) = \lim_{h \downarrow 0} \frac{\mathbb{P}(Y \leq y, x \leq X \leq x+h)}{\mathbb{P}(x \leq X \leq x+h)}.$$

In many statistical situations, the conditional and marginal pdfs are known and (3.9) is used to find the joint pdf via

$$f(x, y) = f_X(x) f_{Y|X}(y|x),$$

or, more generally for the n -dimensional case,

$$f(x_1, \dots, x_n) = f_{X_1}(x_1) f_{X_2|X_1}(x_2|x_1) \cdots f_{X_n|X_1, \dots, X_{n-1}}(x_n|x_1, \dots, x_{n-1}), \quad (3.10)$$

14 which in the discrete case is just a rephrasing of the *product rule* in terms of probability densities. For independent random variables (3.10) reduces to (3.7). Equation (3.10) also shows how one could sequentially generate a random vector $\mathbf{X} = (X_1, \dots, X_n)$ according to a pdf f , provided that it is possible to generate random variables from the successive conditional distributions, as summarized in the following algorithm.

Algorithm 3.1. (Dependent Random Variable Generation).

1. Generate X_1 from pdf f_{X_1} . Set $t = 1$.
2. While $t < n$, given $X_1 = x_1, \dots, X_t = x_t$, generate X_{t+1} from the conditional pdf $f_{X_{t+1}|X_1, \dots, X_t}(x_{t+1} | x_1, \dots, x_t)$ and set $t = t + 1$.
3. Return $\mathbf{X} = (X_1, \dots, X_n)$.

Example 3.8 (Nonuniform Distribution on Triangle). We select a point (X, Y) from the triangle $(0, 0)-(1, 0)-(1, 1)$ in such a way that X has a uniform distribution on $(0, 1)$ and the conditional distribution of Y given $X = x$ is uniform on $(0, x)$. Figure 3.3 shows the result of 1000 independent draws from the joint pdf $f(x, y) = f_X(x) f_{Y|X}(y|x)$, generated via Algorithm 3.1. It is clear that the points are not uniformly distributed over the triangle.

Random variable X has a uniform distribution on $(0, 1)$; hence, its pdf is $f_X(x) = 1$ on $x \in (0, 1)$. For any fixed $x \in (0, 1)$, the conditional distribution of Y given $X = x$ is uniform on the interval $(0, x)$, which means that

$$f_{Y|X}(y|x) = \frac{1}{x}, \quad 0 < y < x.$$

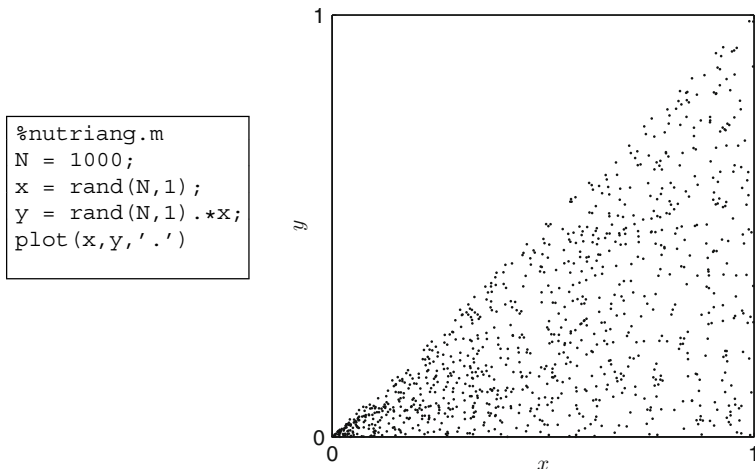


Fig. 3.3 One thousand realizations from the joint density $f(x, y)$, generated using the MATLAB program on the *left*, which implements Algorithm 3.1.

It follows that the joint pdf is given by

$$f(x, y) = f_X(x) f_{Y|X}(y | x) = \frac{1}{x}, \quad 0 < x < 1, \quad 0 < y < x.$$

From the joint pdf we can obtain the pdf of Y as

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx = \int_y^1 \frac{1}{x} dx = -\ln y, \quad 0 < y < 1.$$

Finally, for any fixed $y \in (0, 1)$, the conditional pdf of X given $Y = y$ is

$$f_{X|Y}(x | y) = \frac{f(x, y)}{f_Y(y)} = \frac{-1}{x \ln y}, \quad y < x < 1.$$

3.3 Mixed Joint Distributions

So far we have only considered joint distributions in which the random variables are all discrete or all continuous. The theory can be extended to mixed cases in a straightforward way. For example, the joint pdf of a discrete variable X and a continuous variable Y is defined as the function $f(x, y)$ such that for all events $\{(X, Y) \in A\}$, where $A \subset \mathbb{R}^2$,

$$\mathbb{P}((X, Y) \in A) = \sum_x \int I_{\{(x, y) \in A\}} f(x, y) dy,$$

where I denotes the indicator. The pdf is often specified via (3.10).

Example 3.9 (Beta Distribution). Let $\Theta \sim \mathcal{U}(0, 1)$ and $(X | \Theta = \theta) \sim \text{Bin}(n, \theta)$. Using (3.10), the joint pdf of X and Θ is given by

$$f(x, \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}, \quad \theta \in (0, 1), \quad x = 0, 1, \dots, n.$$

By integrating out θ , we find the pdf of X :

$$f_X(x) = \int_0^1 \binom{n}{x} \theta^x (1 - \theta)^{n-x} d\theta = \binom{n}{x} B(x+1, n-x+1),$$

where B is the **beta function**, defined as

$$B(\alpha, \beta) = \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}, \quad (3.11)$$

48 and Γ is the gamma function in (2.20). The conditional pdf of Θ given $X = x$, where $x \in \{0, \dots, n\}$, is

$$f_{\Theta|X}(\theta | x) = \frac{f(\theta, x)}{f_X(x)} = \frac{\theta^x (1 - \theta)^{n-x}}{B(x+1, n-x+1)}, \quad \theta \in (0, 1).$$

The continuous distribution with pdf

$$f(x; \alpha, \beta) = \frac{\theta^{\alpha-1} (1 - \theta)^{\beta-1}}{B(\alpha, \beta)}, \quad x \in (0, 1) \quad (3.12)$$

is called the **beta distribution** with parameters α and β . Both parameters are assumed to be strictly positive. We write $\text{Beta}(\alpha, \beta)$ for this distribution. For this example we have thus $(\Theta | X = x) \sim \text{Beta}(x+1, n-x+1)$.

3.4 Expectations for Joint Distributions

31 Similar to the univariate case in Theorem 2.2, the expected value of a real-valued function h of $(X_1, \dots, X_n) \sim f$ is a weighted average of all values that $h(X_1, \dots, X_n)$ can take. Specifically, in the continuous case,

$$\mathbb{E}h(X_1, \dots, X_n) = \int \cdots \int h(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \cdots dx_n. \quad (3.13)$$

In the discrete case replace the integrals above with sums.

Two important special cases are the expectation of the *sum* (or more generally affine transformations) of random variables and the *product* of random variables.

Theorem 3.2. (Properties of the Expectation). Let X_1, \dots, X_n be random variables with expectations μ_1, \dots, μ_n . Then,

$$\mathbb{E}[a + b_1 X_1 + b_2 X_2 + \dots + b_n X_n] = a + b_1 \mu_1 + \dots + b_n \mu_n \quad (3.14)$$

for all constants a, b_1, \dots, b_n . Also, for *independent* random variables,

$$\mathbb{E}[X_1 X_2 \dots X_n] = \mu_1 \mu_2 \dots \mu_n. \quad (3.15)$$

Proof. We show it for the continuous case with two variables only. The general case follows by analogy and, for the discrete case, by replacing integrals with sums. Let X_1 and X_2 be continuous random variables with joint pdf f . Then, by (3.13),

$$\begin{aligned} \mathbb{E}[a + b_1 X_1 + b_2 X_2] &= \iint (a + b_1 x_1 + b_2 x_2) f(x_1, x_2) dx_1 dx_2 \\ &= a + b_1 \iint x_1 f(x_1, x_2) dx_1 dx_2 + b_2 \iint x_2 f(x_1, x_2) dx_1 dx_2 \\ &= a + b_1 \int x_1 \left(\int f(x_1, x_2) dx_2 \right) dx_1 + b_2 \int x_2 \left(\int f(x_1, x_2) dx_1 \right) dx_2 \\ &= a + b_1 \int x_1 f_{X_1}(x_1) dx_1 + b_2 \int x_2 f_{X_2}(x_2) dx_2 = a + b_1 \mu_1 + b_2 \mu_2. \end{aligned}$$

Next, assume that X_1 and X_2 are independent, so that $f(x_1, x_2) = f_{X_1}(x_1) \times f_{X_2}(x_2)$. Then,

$$\begin{aligned} \mathbb{E}[X_1 X_2] &= \iint x_1 x_2 f_{X_1}(x_1) f_{X_2}(x_2) dx_1 dx_2 \\ &= \int x_1 f_{X_1}(x_1) dx_1 \times \int x_2 f_{X_2}(x_2) dx_2 = \mu_1 \mu_2. \quad \square \end{aligned}$$

Definition 3.7. (Covariance). The **covariance** of two random variables X and Y with expectations $\mathbb{E}X = \mu_X$ and $\mathbb{E}Y = \mu_Y$ is defined as

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

The covariance is a measure of the amount of linear dependency between two random variables. A scaled version of the covariance is given by the **correlation coefficient**:

$$\varrho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}, \quad (3.16)$$

95 where $\sigma_X^2 = \text{Var}(X)$ and $\sigma_Y^2 = \text{Var}(Y)$. The correlation coefficient always lies between -1 and 1 ; see Problem 3.16.

For easy reference Theorem 3.3 lists some important properties of the variance and covariance.

Theorem 3.3. (Properties of the Variance and Covariance). For random variables X , Y , and Z , and constants a and b , we have

1. $\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2$.
2. $\text{Var}(a + bX) = b^2\text{Var}(X)$.
3. $\text{Cov}(X, Y) = \mathbb{E}XY - \mathbb{E}X \mathbb{E}Y$.
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. If X and Y are independent, then $\text{Cov}(X, Y) = 0$.

33 *Proof.* For simplicity of notation we write $\mathbb{E}Z = \mu_Z$ for a generic random variable Z . Properties 1 and 2 were already shown in Theorem 2.4.

3. $\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}[X Y - X \mu_Y - Y \mu_X + \mu_X \mu_Y] = \mathbb{E}[X Y] - \mu_X \mu_Y$.
4. $\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}[(Y - \mu_Y)(X - \mu_X)] = \text{Cov}(Y, X)$.
5. $\text{Cov}(aX + bY, Z) = \mathbb{E}[(aX + bY)Z] - \mathbb{E}[aX + bY] \mathbb{E}Z = a \mathbb{E}[XZ] - a \mathbb{E}X \mathbb{E}Z + b \mathbb{E}[YZ] - b \mathbb{E}Y \mathbb{E}Z = a \text{Cov}(X, Z) + b \text{Cov}(Y, Z)$.
6. $\text{Cov}(X, X) = \mathbb{E}[(X - \mu_X)(X - \mu_X)] = \mathbb{E}[(X - \mu_X)^2] = \text{Var}(X)$.
7. By Property 6, $\text{Var}(X + Y) = \text{Cov}(X + Y, X + Y)$. By Property 5, $\text{Cov}(X + Y, X + Y) = \text{Cov}(X, X) + \text{Cov}(Y, Y) + \text{Cov}(X, Y) + \text{Cov}(Y, X) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$, where in the last equation Properties 4 and 6 are used.
8. If X and Y are independent, then $\mathbb{E}[X Y] = \mu_X \mu_Y$. Therefore, $\text{Cov}(X, Y) = 0$ follows immediately from Property 3. \square

As a consequence of Properties 2 and 7, we have the following general result for the variance of affine transformations of random variables.

Corollary 3.1. (Variance of an Affine Transformation). Let X_1, \dots, X_n be random variables with variances $\sigma_1^2, \dots, \sigma_n^2$. Then,

$$\text{Var} \left(a + \sum_{i=1}^n b_i X_i \right) = \sum_{i=1}^n b_i^2 \sigma_i^2 + 2 \sum_{i < j} b_i b_j \text{Cov}(X_i, X_j) \quad (3.17)$$

for any choice of constants a and b_1, \dots, b_n . In particular, for *independent* random variables X_1, \dots, X_n ,

$$\text{Var}(a + b_1 X_1 + \dots + b_n X_n) = b_1^2 \sigma_1^2 + \dots + b_n^2 \sigma_n^2. \quad (3.18)$$

Let $\mathbf{X} = (X_1, \dots, X_n)^\top$ be a random column vector. Sometimes it is convenient to write the expectations and covariances in vector notation.

Definition 3.8. (Expectation Vector and Covariance Matrix). For any random column vector \mathbf{X} we define the **expectation vector** as the vector of expectations

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top = (\mathbb{E}X_1, \dots, \mathbb{E}X_n)^\top.$$

The **covariance matrix** Σ is defined as the matrix whose (i, j) th element is

$$\text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)].$$

If we define the expectation of a matrix to be the matrix of expectations, then we can write the covariance matrix succinctly as

$$\Sigma = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top].$$

Definition 3.9. (Conditional Expectation). The **conditional expectation** of Y given $X = x$, denoted $\mathbb{E}[Y | X = x]$, is the expectation corresponding to the conditional pdf $f_{Y|X}(y | x)$. That is, in the continuous case,

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

In the discrete case replace the integral with a sum.

Note that $\mathbb{E}[Y | X = x]$ is a function of x , say $h(x)$. The corresponding random variable $h(X)$ is written as $\mathbb{E}[Y | X]$. The expectation of $\mathbb{E}[Y | X]$ is, in the continuous case,

$$\begin{aligned}\mathbb{E}\mathbb{E}[Y | X] &= \int \mathbb{E}[Y | X = x] f_X(x) dx = \int \int y \frac{f(x, y)}{f_X(x)} f_X(x) dy dx \\ &= \int y f_Y(y) dy = \mathbb{E}Y .\end{aligned}\tag{3.19}$$

This “stacking” of (conditional) expectations is sometimes referred to as the **tower property**.

Example 3.10 (Nonuniform Distribution on Triangle Continued). In Example 3.8 the conditional expectation of Y given $X = x$, with $0 < x < 1$, is

$$\mathbb{E}[Y | X = x] = \frac{1}{2} x ,$$

because conditioned on $X = x$, Y is uniformly distributed on the interval $(0, x)$. Using the tower property we find

$$\mathbb{E}Y = \frac{1}{2} \mathbb{E}X = \frac{1}{4} .$$

3.5 Functions of Random Variables

Suppose X_1, \dots, X_n are measurements of a random experiment. What can be said about the distribution of a *function* of the data, say $Z = g(X_1, \dots, X_n)$, when the joint distribution of X_1, \dots, X_n is known?

Example 3.11 (Pdf of an Affine Transformation). Let X be a continuous random variable with pdf f_X and let $Z = a + bX$, where $b \neq 0$. We wish to determine the pdf f_Z of Z . Suppose that $b > 0$. We have for any z

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

Differentiating this with respect to z gives $f_Z(z) = f_X((z - a)/b) / b$. For $b < 0$ we similarly obtain $f_Z(z) = f_X((z - a)/b) / (-b)$. Thus, in general,

$$f_Z(z) = \frac{1}{|b|} f_X\left(\frac{z - a}{b}\right) .\tag{3.20}$$

Example 3.12 (Pdf of a Monotone Transformation). Generalizing the previous example, suppose that $Z = g(X)$ for some strictly increasing function g . To find the pdf of Z from that of X we first write

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

where g^{-1} is the inverse of g . Differentiating with respect to z now gives

$$f_Z(z) = f_X(g^{-1}(z)) \frac{d}{dz} g^{-1}(z) = \frac{f_X(g^{-1}(z))}{g'(g^{-1}(z))}. \quad (3.21)$$

For strictly decreasing functions, g' needs to be replaced with its negative value.

3.5.1 Linear Transformations

Let $\mathbf{x} = (x_1, \dots, x_n)^\top$ be a column vector in \mathbb{R}^n and B an $m \times n$ matrix. The mapping $\mathbf{x} \mapsto \mathbf{z}$, with $\mathbf{z} = B\mathbf{x}$, is called a **linear transformation**. Now consider a *random* vector $\mathbf{X} = (X_1, \dots, X_n)^\top$, and let

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

Then \mathbf{Z} is a random vector in \mathbb{R}^m . In principle, if we know the joint distribution of \mathbf{X} , then we can derive the joint distribution of \mathbf{Z} . Let us first see how the expectation vector and covariance matrix are transformed.

Theorem 3.4. (Expectation and Covariance Under a Linear Transformation). If \mathbf{X} has expectation vector $\boldsymbol{\mu}_X$ and covariance matrix Σ_X , then the expectation vector and covariance matrix of $\mathbf{Z} = B\mathbf{X}$ are given by

$$\boldsymbol{\mu}_Z = B\boldsymbol{\mu}_X \quad (3.22)$$

and

$$\Sigma_Z = B\Sigma_X B^\top. \quad (3.23)$$

Proof. We have $\mu_Z = \mathbb{E}Z = \mathbb{E}B\mathbf{X} = B \mathbb{E}\mathbf{X} = B\mu_X$ and

$$\begin{aligned}\Sigma_Z &= \mathbb{E}[(Z - \mu_Z)(Z - \mu_Z)^\top] = \mathbb{E}[B(\mathbf{X} - \mu_X)(B(\mathbf{X} - \mu_X))^\top] \\ &= B \mathbb{E}[(\mathbf{X} - \mu_X)(\mathbf{X} - \mu_X)^\top] B^\top \\ &= B \Sigma_X B^\top .\end{aligned}$$

□

Suppose that B is an *invertible* $n \times n$ matrix. If \mathbf{X} has a joint pdf f_X , what is the joint density f_Z of \mathbf{Z} ? Let us consider the continuous case. For any fixed \mathbf{x} , let $\mathbf{z} = B\mathbf{x}$. Hence, $\mathbf{x} = B^{-1}\mathbf{z}$. Consider the n -dimensional cube $C = [z_1, z_1 + h] \times \cdots \times [z_n, z_n + h]$. Then, by definition of the joint density for \mathbf{Z} , we have

$$\mathbb{P}(\mathbf{Z} \in C) \approx h^n f_Z(\mathbf{z}) .$$

Let D be the image of C under B^{-1} —that is, the parallelepiped of all points \mathbf{x} such that $B\mathbf{x} \in C$; see Fig. 3.4.

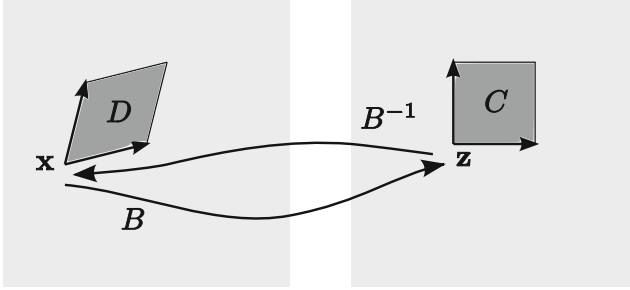


Fig. 3.4 Linear transformation

A basic result from linear algebra is that any matrix B linearly transforms an n -dimensional rectangle with volume V into an n -dimensional parallelepiped with volume $V |B|$, where $|B| = |\det(B)|$. Thus, in addition to the above expression for $\mathbb{P}(\mathbf{Z} \in C)$, we also have

$$\mathbb{P}(\mathbf{Z} \in C) = \mathbb{P}(\mathbf{X} \in D) \approx h^n |B^{-1}| f_X(\mathbf{x}) = h^n |B|^{-1} f_X(\mathbf{x}) .$$

Equating these two expressions for $\mathbb{P}(\mathbf{Z} \in C)$ and letting h go to 0, we obtain

$$f_Z(\mathbf{z}) = \frac{f_X(B^{-1}\mathbf{z})}{|B|}, \quad \mathbf{z} \in \mathbb{R}^n . \quad (3.24)$$

3.5.2 General Transformations

We can apply similar reasoning as in the previous subsection to deal with general transformations $\mathbf{x} \mapsto \mathbf{g}(\mathbf{x})$, written out as

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \mapsto \begin{pmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \\ \vdots \\ g_n(\mathbf{x}) \end{pmatrix} .$$

For a fixed \mathbf{x} , let $\mathbf{z} = \mathbf{g}(\mathbf{x})$. Suppose \mathbf{g} is invertible; hence, $\mathbf{x} = \mathbf{g}^{-1}(\mathbf{z})$. Any infinitesimal n -dimensional rectangle at \mathbf{x} with volume V is transformed into an n -dimensional parallelepiped at \mathbf{z} with volume $V |J_{\mathbf{g}}(\mathbf{x})|$, where $J_{\mathbf{g}}(\mathbf{x})$ is the *matrix of Jacobi* at \mathbf{x} of the transformation \mathbf{g} ; that is,

 367

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

Now consider a random column vector $\mathbf{Z} = \mathbf{g}(\mathbf{X})$. Let C be a small cube around \mathbf{z} with volume h^n . Let D be the image of C under \mathbf{g}^{-1} . Then, as in the linear case,

$$h^n f_{\mathbf{Z}}(\mathbf{z}) \approx \mathbb{P}(\mathbf{Z} \in C) \approx h^n |J_{\mathbf{g}^{-1}}(\mathbf{z})| f_{\mathbf{X}}(\mathbf{x}) .$$

Hence, we have the following result.

Theorem 3.5. (Transformation Rule). Let \mathbf{X} be a continuous n -dimensional random vector with pdf $f_{\mathbf{X}}$ and \mathbf{g} a function from \mathbb{R}^n to \mathbb{R}^n with inverse \mathbf{g}^{-1} . Then, $\mathbf{Z} = \mathbf{g}(\mathbf{X})$ has pdf

$$f_{\mathbf{Z}}(\mathbf{z}) = f_{\mathbf{X}}(\mathbf{g}^{-1}(\mathbf{z})) |J_{\mathbf{g}^{-1}}(\mathbf{z})| , \quad \mathbf{z} \in \mathbb{R}^n . \quad (3.25)$$

Remark 3.3. Note that $|J_{\mathbf{g}^{-1}}(\mathbf{z})| = 1/|J_{\mathbf{g}}(\mathbf{x})|$.

Example 3.13 (Box–Muller Method). The joint distribution of $X, Y \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ is

$$f_{X,Y}(x, y) = \frac{1}{2\pi} e^{-\frac{1}{2}(x^2+y^2)} , \quad (x, y) \in \mathbb{R}^2 .$$

In polar coordinates we have

$$X = R \cos \Theta \quad \text{and} \quad Y = R \sin \Theta , \quad (3.26)$$

where $R \geq 0$ and $\Theta \in (0, 2\pi)$. What is the joint pdf of R and Θ ? Consider the inverse transformation \mathbf{g}^{-1} , defined by

$$\begin{pmatrix} r \\ \theta \end{pmatrix} \xrightarrow{\mathbf{g}^{-1}} \begin{pmatrix} r \cos \theta \\ r \sin \theta \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix}.$$

The corresponding matrix of Jacobi is

$$J_{\mathbf{g}^{-1}}(r, \theta) = \begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix},$$

which has determinant r . Since $x^2 + y^2 = r^2(\cos^2 \theta + \sin^2 \theta) = r^2$, it follows that

$$f_{R,\Theta}(r, \theta) = f_{X,Y}(x, y) r = \frac{1}{2\pi} e^{-\frac{1}{2}r^2} r, \quad \theta \in (0, 2\pi), \quad r \geq 0.$$

By integrating out θ and r , respectively, we find $f_R(r) = r e^{-r^2/2}$ and $f_\Theta(\theta) = 1/(2\pi)$. Since $f_{R,\Theta}$ is the product of f_R and f_Θ , the random variables R and Θ are independent. This shows how X and Y could be generated: independently generate $R \sim f_R$ and $\Theta \sim \mathcal{U}(0, 2\pi)$ and return X and Y via (3.26). Generation from f_R can be done via the inverse-transform method. In particular, R has the same distribution as $\sqrt{-2 \ln U}$ with $U \sim \mathcal{U}(0, 1)$. This leads to the following method for generating standard normal random variables.

53

Algorithm 3.2. (Box–Muller Method).

1. Generate $U_1, U_2 \stackrel{\text{iid}}{\sim} \mathcal{U}(0, 1)$.
2. Return two independent standard normal variables, X and Y , via

$$\begin{aligned} X &= \sqrt{-2 \ln U_1} \cos(2\pi U_2), \\ Y &= \sqrt{-2 \ln U_1} \sin(2\pi U_2). \end{aligned} \tag{3.27}$$

3.6 Multivariate Normal Distribution

It is helpful to view a normally distributed random variable as an affine transformation of a standard normal random variable. In particular, if Z has a standard normal distribution, then $X = \mu + \sigma Z$ has a $\mathcal{N}(\mu, \sigma^2)$ distribution; see Theorem 2.15.

46

We now generalize this to n dimensions. Let Z_1, \dots, Z_n be independent and standard normal random variables. The joint pdf of $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$ is given by

$$f_{\mathbf{Z}}(\mathbf{z}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_i^2} = (2\pi)^{-\frac{n}{2}} e^{-\frac{1}{2}\mathbf{z}^\top \mathbf{z}}, \quad \mathbf{z} \in \mathbb{R}^n. \quad (3.28)$$

We write $\mathbf{Z} \sim \mathbf{N}(\mathbf{0}, I)$, where I is the identity matrix. Consider the affine transformation (i.e., a linear transformation plus a constant vector)

$$\mathbf{X} = \boldsymbol{\mu} + B\mathbf{Z} \quad (3.29)$$

for some $m \times n$ matrix B and m -dimensional vector $\boldsymbol{\mu}$. Note that, by Theorem 3.4, \mathbf{X} has expectation vector $\boldsymbol{\mu}$ and covariance matrix $\Sigma = BB^\top$.

Definition 3.10. (Multivariate Normal Distribution). A random vector \mathbf{X} is said to have a **multivariate normal** or **multivariate Gaussian** distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix Σ if it can be written as $\mathbf{X} = \boldsymbol{\mu} + B\mathbf{Z}$, where $\mathbf{Z} \sim \mathbf{N}(\mathbf{0}, I)$ and $BB^\top = \Sigma$. We write $\mathbf{X} \sim \mathbf{N}(\boldsymbol{\mu}, \Sigma)$.

Suppose that B is an invertible $n \times n$ matrix. Then, by (3.24), the density of $\mathbf{Y} = \mathbf{X} - \boldsymbol{\mu}$ is given by

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

We have $|B| = \sqrt{|\Sigma|}$ and $(\mathbf{B}^{-1})^\top \mathbf{B}^{-1} = (\mathbf{B}^\top)^{-1} \mathbf{B}^{-1} = (\mathbf{B}\mathbf{B}^\top)^{-1} = \Sigma^{-1}$, so that

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

Because \mathbf{X} is obtained from \mathbf{Y} by simply adding a constant vector $\boldsymbol{\mu}$, we have $f_{\mathbf{X}}(\mathbf{x}) = f_{\mathbf{Y}}(\mathbf{x} - \boldsymbol{\mu})$ and therefore

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})}, \quad \mathbf{x} \in \mathbb{R}^n. \quad (3.30)$$

Figure 3.5 shows the pdfs of two bivariate (i.e., two-dimensional) normal distributions. In both cases the mean vector is $\boldsymbol{\mu} = (0, 0)^\top$ and the variances (the diagonal elements of Σ) are 1. The correlation coefficients (or, equivalently here, the covariances) are, respectively, $\varrho = 0$ and $\varrho = 0.8$.

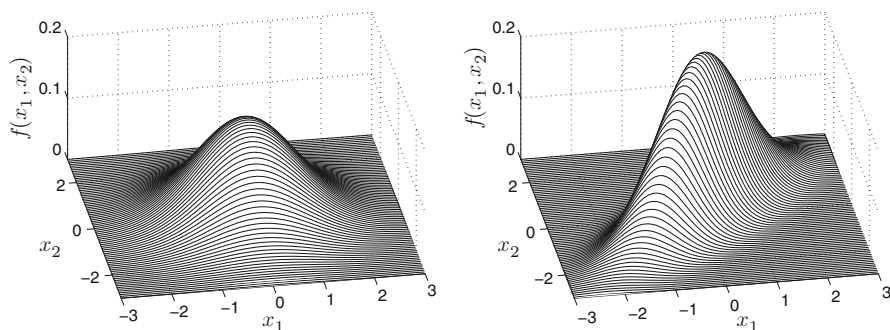


Fig. 3.5 Pdfs of bivariate normal distributions with means zero, variances 1, and correlation coefficients 0 (*left*) and 0.8 (*right*)

Conversely, given a covariance matrix $\Sigma = (\sigma_{ij})$, there exists a unique lower triangular matrix B such that $\Sigma = BB^\top$. In MATLAB, the function `chol` accomplishes this so-called **Cholesky factorization**. Note that it is important to use the option `'lower'` when calling this function, as MATLAB produces an upper triangular matrix by default. Once the Cholesky factorization is determined, it is easy to sample from a multivariate normal distribution.

Algorithm 3.3. (Normal Random Vector Generation). To generate N independent draws from a $N(\mu, \Sigma)$ distribution of dimension n carry out the following steps:

1. Determine the lower Cholesky factorization $\Sigma = BB^\top$.
2. Generate $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$ by drawing $Z_1, \dots, Z_n \sim_{\text{iid}} N(0, 1)$.
3. Output $\mathbf{X} = \mu + B\mathbf{Z}$.
4. Repeat Steps 2 and 3 independently N times.

Example 3.14 (Generating from a Bivariate Normal Distribution). The MATLAB code below draws 1000 samples from the two pdfs in Fig. 3.5. The resulting point clouds are given in Fig. 3.6.

```
%bivnorm.m
N = 1000; rho = 0.8;
Sigma = [1 rho; rho 1];
B=chol(Sigma,'lower');
x=B*randn(2,N);
plot(x(1,:),x(2,:),'.')
```

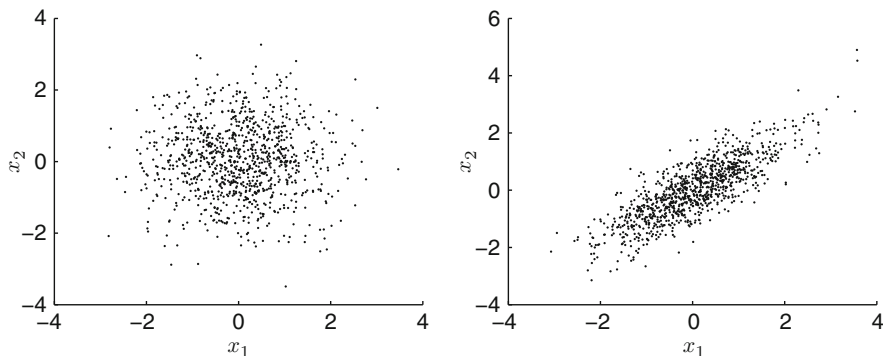


Fig. 3.6 One thousand realizations of bivariate normal distributions with means zero, variances 1, and correlation coefficients 0 (*left*) and 0.8 (*right*)

The following theorem states that any affine combination of independent multivariate normal random variables is again multivariate normal.

Theorem 3.6. (Affine Transformation of Normal Random Vectors). Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_r$ be independent m_i -dimensional normal random vectors, with $\mathbf{X}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \Sigma_i)$, $i = 1, \dots, r$. Then, for any $n \times 1$ vector \mathbf{a} and $n \times m_i$ matrices B_1, \dots, B_r ,

$$\mathbf{a} + \sum_{i=1}^r B_i \mathbf{X}_i \sim \mathcal{N}\left(\mathbf{a} + \sum_{i=1}^r B_i \boldsymbol{\mu}_i, \sum_{i=1}^r B_i \Sigma_i B_i^\top\right). \quad (3.31)$$

Proof. Denote the n -dimensional random vector in the left-hand side of (3.31) by \mathbf{Y} . By Definition 3.10, each \mathbf{X}_i can be written as $\boldsymbol{\mu}_i + A_i \mathbf{Z}_i$, where the $\{\mathbf{Z}_i\}$ are independent (because the $\{\mathbf{X}_i\}$ are independent), so that

$$\mathbf{Y} = \mathbf{a} + \sum_{i=1}^r B_i (\boldsymbol{\mu}_i + A_i \mathbf{Z}_i) = \mathbf{a} + \sum_{i=1}^r B_i \boldsymbol{\mu}_i + \sum_{i=1}^r B_i A_i \mathbf{Z}_i,$$

which is an affine combination of independent standard normal random vectors. Hence, \mathbf{Y} is multivariate normal. Its expectation vector and covariance matrix can be found easily from Theorem 3.4. \square

The next theorem shows that the distribution of a subvector of a multivariate normal random vector is again normal.

Theorem 3.7. (Marginal Distributions of Normal Random Vectors). Let $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ be an n -dimensional normal random vector. Decompose \mathbf{X} , $\boldsymbol{\mu}$, and Σ as

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_p \\ \mathbf{X}_q \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_p \\ \boldsymbol{\mu}_q \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_p & \Sigma_r \\ \Sigma_r^\top & \Sigma_q \end{pmatrix}, \quad (3.32)$$

where Σ_p is the upper left $p \times p$ corner of Σ and Σ_q is the lower right $q \times q$ corner of Σ . Then, $\mathbf{X}_p \sim \mathcal{N}(\boldsymbol{\mu}_p, \Sigma_p)$.

Proof. Let BB^\top be the lower Cholesky factorization of Σ . We can write

$$\begin{pmatrix} \mathbf{X}_p \\ \mathbf{X}_q \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_p \\ \boldsymbol{\mu}_q \end{pmatrix} + \underbrace{\begin{pmatrix} B_p & O \\ C_r & C_q \end{pmatrix}}_B \begin{pmatrix} \mathbf{Z}_p \\ \mathbf{Z}_q \end{pmatrix}, \quad (3.33)$$

where \mathbf{Z}_p and \mathbf{Z}_q are independent p - and q -dimensional standard normal random vectors. In particular, $\mathbf{X}_p = \boldsymbol{\mu}_p + B_p \mathbf{Z}_p$, which means that $\mathbf{X}_p \sim \mathcal{N}(\boldsymbol{\mu}_p, \Sigma_p)$, since $B_p B_p^\top = \Sigma_p$. \square

By relabeling the elements of \mathbf{X} we see that Theorem 3.7 implies that *any* subvector of \mathbf{X} has a multivariate normal distribution. For example, $\mathbf{X}_q \sim \mathcal{N}(\boldsymbol{\mu}_q, \Sigma_q)$.

Not only the marginal distributions of a normal random vector are normal but also its *conditional distributions*.

Theorem 3.8. (Conditional Distributions of Normal Random Vectors). Let $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ be an n -dimensional normal random vector with $\det(\Sigma) > 0$. If \mathbf{X} is decomposed as in (3.32), then

$$(\mathbf{X}_q | \mathbf{X}_p = \mathbf{x}_p) \sim \mathcal{N}(\boldsymbol{\mu}_q + \Sigma_r^\top \Sigma_p^{-1}(\mathbf{x}_p - \boldsymbol{\mu}_p), \Sigma_q - \Sigma_r^\top \Sigma_p^{-1} \Sigma_r). \quad (3.34)$$

As a consequence, \mathbf{X}_p and \mathbf{X}_q are *independent* if and only if they are *uncorrelated*; that is, if $\Sigma_r = O$ (zero matrix).

Proof. From (3.33) we see that

$$(\mathbf{X}_q | \mathbf{X}_p = \mathbf{x}_p) = \boldsymbol{\mu}_q + C_r B_p^{-1}(\mathbf{x}_p - \boldsymbol{\mu}_p) + C_q \mathbf{Z}_q,$$

where \mathbf{Z}_q is a q -dimensional multivariate standard normal random vector. It follows that \mathbf{X}_q conditional on $\mathbf{X}_p = \mathbf{x}_p$ has a $\mathcal{N}(\boldsymbol{\mu}_q + C_r B_p^{-1}(\mathbf{x}_p - \boldsymbol{\mu}_p), C_q C_q^\top)$ distribution.

The proof of (3.34) is completed by observing that $\Sigma_r^\top \Sigma_p^{-1} = C_r B_p^\top (B_p^\top)^{-1} B_p^{-1} = C_r B_p^{-1}$, and

$$\Sigma_q - \Sigma_r^\top \Sigma_p^{-1} \Sigma_r = C_r C_r^\top + C_q C_q^\top - C_r B_p^{-1} \underbrace{\Sigma_r}_{B_p C_r^\top} = C_q C_q^\top.$$

If \mathbf{X}_p and \mathbf{X}_q are independent, then they are obviously uncorrelated, as $\Sigma_r = \mathbb{E}[(\mathbf{X}_p - \boldsymbol{\mu}_p)(\mathbf{X}_q - \boldsymbol{\mu}_q)^\top] = \mathbb{E}(\mathbf{X}_p - \boldsymbol{\mu}_p) \mathbb{E}(\mathbf{X}_q - \boldsymbol{\mu}_q)^\top = \mathbf{O}$. Conversely, if $\Sigma_r = \mathbf{O}$, then by (3.34), the conditional distribution of \mathbf{X}_q given \mathbf{X}_p is the same as the unconditional distribution of \mathbf{X}_q , that is, $N(\boldsymbol{\mu}_q, \Sigma_q)$. In other words, \mathbf{X}_q is independent of \mathbf{X}_p . \square

Theorem 3.9. (Relationship Between Normal and χ^2 Distributions). If $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$ is an n -dimensional normal random with vector with $\det(\Sigma) > 0$, then

$$(\mathbf{X} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{X} - \boldsymbol{\mu}) \sim \chi_n^2. \quad (3.35)$$

Proof. Let BB^\top be the Cholesky factorization of Σ , where B is invertible. Since \mathbf{X} can be written as $\boldsymbol{\mu} + B\mathbf{Z}$, where $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$ is a vector of independent standard normal random variables, we have

$$(\mathbf{X} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{X} - \boldsymbol{\mu}) = (\mathbf{X} - \boldsymbol{\mu})^\top (BB^\top)^{-1} (\mathbf{X} - \boldsymbol{\mu}) = \mathbf{Z}^\top \mathbf{Z} = \sum_{i=1}^n Z_i^2.$$

The moment generating function of $Y = \sum_{i=1}^n Z_i^2$ is given by

$$\mathbb{E} e^{tY} = \mathbb{E} e^{t(Z_1^2 + \dots + Z_n^2)} = \mathbb{E} [e^{tZ_1^2} \dots e^{tZ_n^2}] = \left(\mathbb{E} e^{tZ^2} \right)^n,$$

where $Z \sim N(0, 1)$. The moment generating function of Z^2 is

$$\mathbb{E} e^{tZ^2} = \int_{-\infty}^{\infty} e^{tz^2} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(1-2t)z^2} dz = \frac{1}{\sqrt{1-2t}},$$

so that

$$\mathbb{E} e^{tY} = \left(\frac{\frac{1}{2}}{\frac{1}{2} - t} \right)^{\frac{n}{2}}, \quad t < \frac{1}{2},$$

which is the moment generating function of the Gamma($n/2, 1/2$) distribution, that is, the χ_n^2 distribution—see Theorem 2.18. \square

A consequence of Theorem 3.9 is that if $\mathbf{X} = (X_1, \dots, X_n)^\top$ is n -dimensional standard normal, then the squared length $\|\mathbf{X}\|^2 = X_1^2 + \dots + X_n^2$ has a χ_n^2 distribution. If instead $X_i \sim N(\mu_i, 1)$, $i = 1, \dots, n$, then $\|\mathbf{X}\|^2$ is said to have a **noncentral χ_n^2 distribution**. This distribution depends on the $\{\mu_i\}$ only through the norm $\|\boldsymbol{\mu}\|$; see Problem 3.22. We write $\|\mathbf{X}\|^2 \sim \chi_n^2(\theta)$, where $\theta = \|\boldsymbol{\mu}\|^2$ is the **noncentrality parameter**.

Such distributions frequently occur in statistics when considering *projections* of multivariate normal random variables. The proof of the following theorem can be found in Appendix B.4.

71

Theorem 3.10. (Relationship Between Normal and Noncentral χ^2 Distributions). Let $\mathbf{X} \sim N(\boldsymbol{\mu}, I)$ be an n -dimensional normal random vector and let \mathcal{V}_k and \mathcal{V}_m be linear subspaces of dimensions k and m , respectively, with $k < m \leq n$. Let \mathbf{X}_k and \mathbf{X}_m be orthogonal projections of \mathbf{X} onto \mathcal{V}_k and \mathcal{V}_m , and let $\boldsymbol{\mu}_k$ and $\boldsymbol{\mu}_m$ be the corresponding projections of $\boldsymbol{\mu}$. Then, the following holds:

1. The random vectors \mathbf{X}_k , $\mathbf{X}_m - \mathbf{X}_k$, and $\mathbf{X} - \mathbf{X}_m$ are independent.
2. $\|\mathbf{X}_k\|^2 \sim \chi_k^2(\|\boldsymbol{\mu}_k\|)$, $\|\mathbf{X}_m - \mathbf{X}_k\|^2 \sim \chi_{m-k}^2(\|\boldsymbol{\mu}_m - \boldsymbol{\mu}_k\|)$, and $\|\mathbf{X} - \mathbf{X}_m\|^2 \sim \chi_{n-m}^2(\|\boldsymbol{\mu} - \boldsymbol{\mu}_m\|)$.

Theorem 3.10 is frequently used in the statistical analysis of *normal linear models*; see Sect. 5.3.1. In typical situations $\boldsymbol{\mu}$ lies in the subspace \mathcal{V}_m or even \mathcal{V}_k —in which case $\|\mathbf{X}_m - \mathbf{X}_k\|^2 \sim \chi_{m-k}^2$ and $\|\mathbf{X} - \mathbf{X}_m\|^2 \sim \chi_{n-m}^2$, independently. The (scaled) quotient then turns out to have an F distribution—a consequence of the following theorem.

142

Theorem 3.11. (Relationship Between χ^2 and F Distributions). Let $U \sim \chi_m^2$ and $V \sim \chi_n^2$ be independent. Then,

$$\frac{U/m}{V/n} \sim F(m, n).$$

Proof. For notational simplicity, let $c = m/2$ and $d = n/2$. It follows from Example 3.7 that the pdf of $W = U/V$ is given by

71

$$f_W(w) = \int_0^\infty f_U(wv) v f_V(v) dv$$

$$\begin{aligned}
&= \int_0^\infty \frac{(wv)^{c-1} e^{-wv/2}}{\Gamma(c) 2^c} v \frac{v^{d-1} e^{-v/2}}{\Gamma(d) 2^d} dv \\
&= \frac{w^{c-1}}{\Gamma(c) \Gamma(d) 2^{c+d}} \int_0^\infty v^{c+d-1} e^{-(1+w)v/2} dv \\
&= \frac{\Gamma(c+d)}{\Gamma(c) \Gamma(d)} \frac{w^{c-1}}{(1+w)^{c+d}},
\end{aligned}$$

where the last equality follows from the fact that the integrand is equal to $\Gamma(\alpha)\lambda^\alpha$ times the density of the Gamma(α, λ) distribution with $\alpha = c + d$ and $\lambda = (1 + w)/2$. The proof is completed by observing that the density of $Z = \frac{n}{m} \frac{U}{V}$ is given by

$$f_Z(z) = f_W(zm/n) m/n.$$

□

Corollary 3.2. (Relationship Between Normal, χ^2 , and t Distributions).

Let $Z \sim N(0, 1)$ and $V \sim \chi_n^2$ be independent. Then,

$$\frac{Z}{\sqrt{V/n}} \sim t_n.$$

Proof. Let $T = Z/\sqrt{V/n}$. Because $Z^2 \sim \chi_1^2$, we have by Theorem 3.11 that $T^2 \sim F(1, n)$. The result follows now from Theorem 2.19 and the symmetry around 0 of the pdf of T . □

51

3.7 Limit Theorems

Two main results in probability are the *law of large numbers* and the *central limit theorem*. Both are limit theorems involving sums of independent random variables. In particular, consider a sequence X_1, X_2, \dots of iid random variables with finite expectation μ and finite variance σ^2 . For each n define $S_n = X_1 + \dots + X_n$. What can we say about the (random) sequence of sums S_1, S_2, \dots or averages $S_1, S_2/2, S_3/3, \dots$? By (3.14) and (3.18) we have $\mathbb{E}[S_n/n] = \mu$ and $\text{Var}(S_n/n) = \sigma^2/n$. Hence, as n increases, the variance of the (random) average S_n/n goes to 0. Informally, this means that (S_n/n) tends to the constant μ , as $n \rightarrow \infty$. This makes intuitive sense, but the important point is that the mathematical theory *confirms* our intuition in this respect. Here is a more precise statement.

75

Theorem 3.12. (Weak Law of Large Numbers). If X_1, \dots, X_n are iid with finite expectation μ and finite variance σ^2 , then for all $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|S_n/n - \mu| > \varepsilon) = 0.$$

Proof. Let $Y = (S_n/n - \mu)^2$ and $\delta = \varepsilon^2$. We have

$$\begin{aligned} \text{Var}(S_n/n) &= \mathbb{E}Y = \mathbb{E}[Y \mathbf{I}_{\{Y > \delta\}}] + \mathbb{E}[Y \mathbf{I}_{\{Y \leq \delta\}}] \geq \mathbb{E}[\delta \mathbf{I}_{\{Y > \delta\}}] + 0 \\ &= \delta \mathbb{P}(Y > \delta) = \varepsilon^2 \mathbb{P}(|S_n/n - \mu| > \varepsilon). \end{aligned}$$

Rearranging gives

$$\mathbb{P}(|S_n/n - \mu| > \varepsilon) \leq \frac{\text{Var}(S_n/n)}{\varepsilon^2} = \frac{\sigma^2}{n \varepsilon^2}.$$

The proof is concluded by observing that $\sigma^2/(n\varepsilon^2)$ goes to 0 as $n \rightarrow \infty$. \square

Remark 3.4. In Theorem 3.12 the qualifier “weak” is used to distinguish the result from the *strong* law of large numbers, which states that

$$\mathbb{P}(\lim_{n \rightarrow \infty} S_n/n = \mu) = 1.$$

In terms of a computer simulation this means that the probability of drawing a sequence for which the sequence of averages fails to converge to μ is zero. The strong law implies the weak law, but is more difficult to prove in its full generality; see, for example, (Feller 1970).

The central limit theorem describes the approximate distribution of S_n (or S_n/n), and it applies to both continuous and discrete random variables. Loosely, it states that

the sum of a large number of iid random variables approximately has a normal distribution.

Specifically, the random variable S_n has a distribution that is approximately normal, with expectation $n\mu$ and variance $n\sigma^2$. A more precise statement is given next.

Theorem 3.13. (Central Limit Theorem). If X_1, \dots, X_n are iid with finite expectation μ and finite variance σ^2 , then for all $x \in \mathbb{R}$,

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

where Φ is the cdf of the standard normal distribution.

Proof. (Sketch) A full proof is out of the scope of this book. However, the main ideas are not difficult. Without loss of generality assume $\mu = 0$ and $\sigma = 1$. This amounts to replacing X_n by $(X_n - \mu)/\sigma$. We also assume, for simplicity, that the moment generating function of X_i is finite in an open interval containing 0, so that we can use Theorem 2.7. We wish to show that the cdf of S_n/\sqrt{n} converges to that of the standard normal distribution. It can be proved (and makes intuitive sense) that this is equivalent (up to some technical conditions) to demonstrating that the corresponding moment generating functions converge. That is, we wish to show that

$$\lim_{n \rightarrow \infty} \mathbb{E} \exp \left(t \frac{S_n}{\sqrt{n}} \right) = e^{\frac{1}{2}t^2}, \quad t \in \mathbb{R},$$

where the right-hand side is the moment generating function of the standard normal distribution. Because $\mathbb{E}X_1 = 0$ and $\mathbb{E}X_1^2 = \text{Var}(X_1) = 1$, we have by Theorem 2.7 that the moment generation function of X_1 has the following Taylor expansion:

$$M(t) \stackrel{\text{def}}{=} \mathbb{E} e^{tX_1} = 1 + t \mathbb{E}X_1 + \frac{1}{2}t^2 \mathbb{E}X_1^2 + o(t^2) = 1 + \frac{1}{2}t^2 + o(t^2),$$

where $o(t^2)$ is a function for which $\lim_{t \downarrow 0} o(t^2)/t^2 = 0$. Because the $\{X_i\}$ are iid, it follows that the moment generating function of S_n/\sqrt{n} satisfies

$$\begin{aligned} \mathbb{E} \exp \left(t \frac{S_n}{\sqrt{n}} \right) &= \mathbb{E} \exp \left(\frac{t}{\sqrt{n}} (X_1 + \cdots + X_n) \right) = \prod_{i=1}^n \mathbb{E} \exp \left(\frac{t}{\sqrt{n}} X_i \right) \\ &= M^n \left(\frac{t}{\sqrt{n}} \right) = \left[1 + \frac{t^2}{2n} + o(t^2/n) \right]^n \rightarrow e^{\frac{1}{2}t^2} \end{aligned}$$

as $n \rightarrow \infty$. □

Figure 3.7 shows central limit theorem in action. The left part shows the pdfs of S_1, \dots, S_4 for the case where the $\{X_i\}$ have a $U[0, 1]$ distribution. The right part shows the same for the $\text{Exp}(1)$ distribution. We clearly see convergence to a bell-shaped curve, characteristic of the normal distribution.

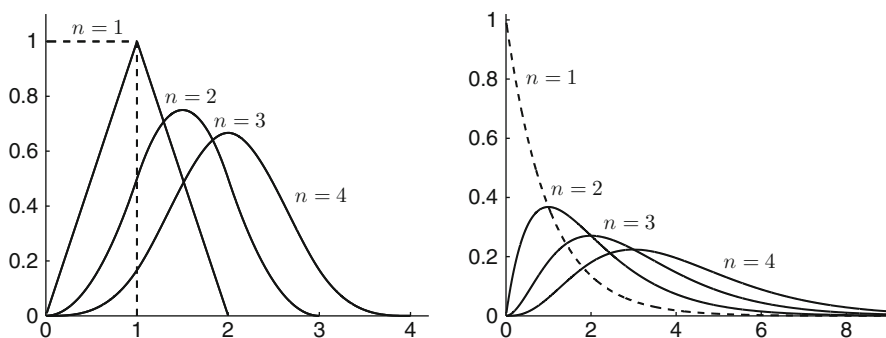


Fig. 3.7 Illustration of the central limit theorem for (left) the uniform distribution and (right) the exponential distribution

66

Recall that a binomial random variable $X \sim \text{Bin}(n, p)$ can be viewed as the sum of n iid $\text{Ber}(p)$ random variables: $X = X_1 + \cdots + X_n$. As a direct consequence of the central limit theorem it follows that, for large n , $\mathbb{P}(X \leq k) \approx \mathbb{P}(Y \leq k)$, where $Y \sim \text{N}(np, np(1-p))$. As a rule of thumb, this normal approximation to the binomial distribution is accurate if both np and $n(1-p)$ are larger than 5.

There is also a central limit theorem for random vectors. The multidimensional version is as follows.

Theorem 3.14. (Multivariate Central Limit Theorem). Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be iid random vectors with expectation vector $\boldsymbol{\mu}$ and covariance matrix Σ . For large n the random vector $\mathbf{X}_1 + \cdots + \mathbf{X}_n$ approximately has a $\text{N}(n\boldsymbol{\mu}, n\Sigma)$ distribution.

A more precise formulation of the above theorem is that the average random vector $\mathbf{Z}_n = (\mathbf{X}_1 + \cdots + \mathbf{X}_n)/n$, when rescaled via $\sqrt{n}(\mathbf{Z}_n - \boldsymbol{\mu})$, converges in distribution to a random vector $\mathbf{K} \sim \text{N}(\mathbf{0}, \Sigma)$ as $n \rightarrow \infty$. A useful consequence of this is given next.

Theorem 3.15. (Delta Method). Let $\mathbf{Z}_1, \mathbf{Z}_2, \dots$ be a sequence of random vectors such that $\sqrt{n}(\mathbf{Z}_n - \boldsymbol{\mu}) \rightarrow \mathbf{K} \sim \text{N}(\mathbf{0}, \Sigma)$ as $n \rightarrow \infty$. Then, for any continuously differentiable function \mathbf{g} of \mathbf{Z}_n ,

$$\sqrt{n}(\mathbf{g}(\mathbf{Z}_n) - \mathbf{g}(\boldsymbol{\mu})) \rightarrow \mathbf{R} \sim \text{N}(\mathbf{0}, J \Sigma J^\top), \quad (3.36)$$

where $J = J(\boldsymbol{\mu}) = (\partial g_i(\boldsymbol{\mu})/\partial x_j)$ is the Jacobian matrix of \mathbf{g} evaluated at $\boldsymbol{\mu}$.

369

Proof. (Sketch) A formal proof requires some deeper knowledge of statistical convergence, but the idea of the proof is quite straightforward. The key step is to construct the first-order Taylor expansion (see Theorem B.1) of \mathbf{g} around $\boldsymbol{\mu}$, which yields

$$\mathbf{g}(\mathbf{Z}_n) = \mathbf{g}(\boldsymbol{\mu}) + J(\boldsymbol{\mu})(\mathbf{Z}_n - \boldsymbol{\mu}) + \mathcal{O}(\|\mathbf{Z}_n - \boldsymbol{\mu}\|^2).$$

80

As $n \rightarrow \infty$, the remainder term goes to 0, because $\mathbf{Z}_n \rightarrow \boldsymbol{\mu}$. Hence, the left-hand side of (3.36) is approximately $J \sqrt{n}(\mathbf{Z}_n - \boldsymbol{\mu})$. For large n this converges to a random vector $\mathbf{R} = J \mathbf{K}$, where $\mathbf{K} \sim \text{N}(\mathbf{0}, \Sigma)$. Finally, by Theorem 3.4, we have $\mathbf{R} \sim \text{N}(\mathbf{0}, J \Sigma J^\top)$. \square

Example 3.15 (Ratio Estimator). Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be iid copies of a random vector (X, Y) with mean vector (μ_X, μ_Y) and covariance matrix Σ . Denoting the average of the $\{X_i\}$ and $\{Y_i\}$ by \bar{X} and \bar{Y} , respectively, what can we say about the distribution of \bar{X}/\bar{Y} for large n ?

Let $\mathbf{Z}_n = (\bar{X}, \bar{Y})$ and $\boldsymbol{\mu} = (\mu_X, \mu_Y)$. By the multivariate central limit theorem \mathbf{Z}_n has approximately a $N(\boldsymbol{\mu}, \Sigma/n)$ distribution. More precisely, $\sqrt{n}(\mathbf{Z}_n - \boldsymbol{\mu})$ converges to a $N(\mathbf{0}, \Sigma)$ -distributed random vector.


We apply the delta method using the function $g(x, y) = x/y$, whose Jacobian matrix is

$$J(x, y) = \left(\frac{\partial g(x, y)}{\partial x}, \quad \frac{\partial g(x, y)}{\partial y} \right) = \left(\frac{1}{y}, \quad -\frac{x}{y^2} \right).$$

It follows from (3.36) that $g(\bar{X}, \bar{Y}) = \bar{X}/\bar{Y}$ has approximately a normal distribution with expectation $g(\boldsymbol{\mu}) = \mu_X/\mu_Y$ and variance σ^2/n , where

$$\begin{aligned} \sigma^2 &= J(\boldsymbol{\mu}) \Sigma J^\top(\boldsymbol{\mu}) = \begin{pmatrix} \frac{1}{\mu_Y} & -\frac{\mu_X}{\mu_Y^2} \end{pmatrix} \begin{pmatrix} \text{Var}(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & \text{Var}(Y) \end{pmatrix} \begin{pmatrix} \frac{1}{\mu_Y} \\ -\frac{\mu_X}{\mu_Y^2} \end{pmatrix} \\ &= \left(\frac{\mu_X}{\mu_Y} \right)^2 \left(\frac{\text{Var}(X)}{\mu_X^2} + \frac{\text{Var}(Y)}{\mu_Y^2} - 2 \frac{\text{Cov}(X, Y)}{\mu_X \mu_Y} \right). \end{aligned} \quad (3.37)$$

3.8 Problems

3.1. Let U and V be independent random variables with $\mathbb{P}(U = 1) = \mathbb{P}(V = 1) = 1/4$ and $\mathbb{P}(U = -1) = \mathbb{P}(V = -1) = 3/4$. Define $X = U/V$ and $Y = U + V$. Give the joint discrete pdf of X and Y in table form, as in Table 3.1. Are X and Y independent?  64


3.2. Let $X_1, \dots, X_4 \sim_{\text{iid}} \text{Ber}(p)$.


- Give the joint discrete pdf of X_1, \dots, X_4 .
- Give the joint discrete pdf of X_1, \dots, X_4 given $X_1 + \dots + X_4 = 2$.

3.3. Three identical-looking urns each have 4 balls. Urn 1 has 1 red and 3 white balls, urn 2 has 2 red and 2 white balls, and urn 3 has 3 red and 1 white ball. We randomly select an urn with equal probability. Let X be the number of the urn. We then draw 2 balls from the selected urn. Let Y be the number of red balls drawn. Find the following discrete pdfs:

- The pdf of X .
- The conditional pdf of Y given $X = x$ for $x = 1, 2, 3$.
- The joint pdf of X and Y .

- (d) The pdf of Y .
 (e) The conditional pdf of X given $Y = y$ for $y = 0, 1, 2$.

 **67** **3.4.** We randomly select a point (X, Y) from the triangle $\{(x, y) : x, y \in \{1, \dots, 6\}, y \leq x\}$ (see Fig. 3.1) in the following *nonuniform* way. First, select X discrete uniformly from $\{1, \dots, 6\}$. Then, given $X = x$, select Y discrete uniformly from $\{1, \dots, x\}$. Find the conditional distribution of X given $Y = 1$ and its corresponding conditional expectation.

 **72** **3.5.** We randomly and uniformly select a continuous random vector (X, Y) in the triangle $(0, 0)-(1, 0)-(1, 1)$, the same triangle as in Example 3.8.

- (a) Give the joint pdf of X and Y .
 (b) Calculate the pdf of Y and sketch its graph.
 (c) Specify the conditional pdf of Y given $X = x$ for any fixed $x \in (0, 1)$.
 (d) Determine $\mathbb{E}[Y | X = 1/2]$.

3.6. Let $X \sim U[0, 1]$ and $Y \sim \text{Exp}(1)$ be independent.

- (a) Determine the joint pdf of X and Y and sketch its graph.
 (b) Calculate $\mathbb{P}((X, Y) \in [0, 1] \times [0, 1])$.
 (c) Calculate $\mathbb{P}(X + Y < 1)$.


3.7. Let $X \sim \text{Exp}(\lambda)$ and $Y \sim \text{Exp}(\mu)$ be independent.

- (a) Show that $\min(X, Y)$ also has an exponential distribution, and determine its corresponding parameter.
 (b) Show that

$$\mathbb{P}(X < Y) = \frac{\lambda}{\lambda + \mu}.$$

3.8. Let $X \sim \text{Exp}(1)$ and $(Y | X = x) \sim \text{Exp}(x)$.

- (a) What is the joint pdf of X and Y ?
 (b) What is the marginal pdf of Y ?

 **71** **3.9.** Let $X \sim U(-\pi/2, \pi/2)$. Show that $Y = \tan(X)$ has a Cauchy distribution.

3.10. Let $X \sim \text{Exp}(3)$ and $Y = \ln(X)$. What is the pdf of Y ?

3.11. We draw n numbers independently and uniformly from the interval $[0, 1]$ and denote their sum S_n .

- (a) Determine the pdf of S_2 and sketch its graph.
 (b) What is approximately the distribution of S_{20} ?
 (c) Approximate the probability that the average of the 20 numbers is greater than 0.6.

3.12. A certain type of electrical component has an exponential lifetime distribution with an expected lifetime of $1/2$ year. When the component fails it is immediately

replaced by a second (new) component; when the second component fails, it is replaced by a third, etc. Suppose there are ten such identical components. Let T be the time that the last of the components fails.

- (a) What is the expectation and variance of T ?
- (b) Approximate, using the central limit theorem, the probability that T exceeds 6 years.
- (c) What is the exact distribution of T ?

3.13. Let A be an invertible $n \times n$ matrix and let $X_1, \dots, X_n \sim_{\text{iid}} N(0, 1)$. Define $\mathbf{X} = (X_1, \dots, X_n)^\top$ and let $(Z_1, \dots, Z_n)^\top = A\mathbf{X}$. Show that Z_1, \dots, Z_n are iid standard normal only if $AA^\top = I$ (identity matrix), in other words, only if A is an *orthogonal* matrix. Can you find a geometric interpretation of this?

3.14. Let X_1, \dots, X_n be independent and identically distributed random variables with mean μ and variance σ^2 . Let $\bar{X} = (X_1 + \dots + X_n)/n$. Calculate the correlation coefficient of X_1 and \bar{X} .

3.15. Suppose that X_1, \dots, X_6 are iid with pdf

$$f(x) = \begin{cases} 3x^2, & 0 \leq x \leq 1, \\ 0, & \text{elsewhere.} \end{cases}$$

- (a) What is the probability that all $\{X_i\}$ are greater than $1/2$?
- (b) Find the probability that at least one of the $\{X_i\}$ is less than $1/2$.

3.16. Let X and Y be random variables.

- (a) Express $\text{Var}(-aX + Y)$, where a is a constant, in terms of $\text{Var}(X)$, $\text{Var}(Y)$, and $\text{Cov}(X, Y)$.
- (b) Take $a = \text{Cov}(X, Y)/\text{Var}(X)$. Using the fact that the variance in (a) is always nonnegative, prove the following **Cauchy–Schwarz inequality**:

$$(\text{Cov}(X, Y))^2 \leq \text{Var}(X) \text{Var}(Y).$$

- (c) Show that, as a consequence, the correlation coefficient of X and Y must lie between -1 and 1 .

3.17. Suppose X and Y are independent uniform random variables on $[0, 1]$. Let $U = X/Y$ and $V = XY$, which means $X = \sqrt{UV}$ and $Y = \sqrt{V/U}$.

- (a) Sketch the two-dimensional region where the density of (U, V) is nonzero.
- (b) Find the matrix of Jacobi for the transformation $(x, y)^\top \mapsto (u, v)^\top$.
- (c) Show that its determinant is $2x/y = 2u$.
- (d) What is the joint pdf of U and V ?
- (e) Show that the marginal pdf of U is

$$f_U(u) = \begin{cases} \frac{1}{2}, & 0 < u < 1 \\ \frac{1}{2u^2}, & u \geq 1 \end{cases}. \quad (3.38)$$

3.18. Let X_1, \dots, X_n be iid with mean μ and variance σ^2 . Let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and $Y = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$.

(a) Show that

$$Y = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2.$$

(b) Calculate $\mathbb{E}Y$.

(c) Show that $\mathbb{E}Y \rightarrow \sigma^2$ as $n \rightarrow \infty$.

3.19. Let $\mathbf{X} = (X_1, \dots, X_n)^\top$, with $\{X_i\} \sim_{\text{iid}} N(\mu, 1)$. Consider the orthogonal projection, denoted \mathbf{X}_1 , of \mathbf{X} onto the subspace spanned by $\mathbf{1} = (1, \dots, 1)^\top$.

(a) Show that $\mathbf{X}_1 = \bar{X}\mathbf{1}$.

(b) Show that \mathbf{X}_1 and $\mathbf{X} - \mathbf{X}_1$ are independent.

(c) Show that $\|\mathbf{X} - \mathbf{X}_1\|^2 = \sum_{i=1}^n (X_i - \bar{X})^2$ has a χ_{n-1}^2 distribution.

Hint: apply Theorem 3.10.

3.20. Let X_1, \dots, X_6 be the weights of six randomly chosen people. Assume each weight is $N(75|100)$ distributed (in kg). Let $W = X_1 + \dots + X_6$ be the total weight of the group. Explain why the distribution of W is equal or not equal to $6X_1$.

3.21. Let $X \sim \chi_m^2$ and $Y \sim \chi_n^2$ be independent. Show that $X + Y \sim \chi_{m+n}^2$. Hint: use moment generating functions.

3.22. Let $X \sim N(\mu, 1)$. Show that the moment generation function of X^2 is

$$M(t) = \frac{e^{\mu^2 t / (1-2t)}}{\sqrt{1-2t}} \quad t < 1/2.$$

Next, consider independent random variables $X_i \sim N(\mu_i, 1)$, $i = 1, \dots, n$. Use the result above to show that the distribution of $\|\mathbf{X}\|^2$ only depends on n and $\|\boldsymbol{\mu}\|$. Can you find a symmetry argument why this must be so?

3.23. A machine produces cylinders with a diameter which is normally distributed with mean 3.97 cm and standard deviation 0.03 cm. Another machine produces (independently of the first machine) shafts with a diameter which is normally distributed with mean 4.05 cm and standard deviation 0.02 cm. What is the probability that a randomly chosen cylinder fits into a randomly chosen shaft?

3.24. A sieve with diameter d is used to separate a large number of blueberries into two classes: small and large. Suppose that the diameters of the blueberries are normally distributed with an expectation $\mu = 1$ cm and a standard deviation $\sigma = 0.1$ cm.

(a) Find the diameter of the sieve such that the proportion of large blueberries is 30%.

- (b) Suppose that the diameter is chosen such as in (a). What is the probability that out of 1000 blueberries, fewer than 280 end up in the “large” class?

3.25. Suppose X , Y , and Z are independent $N(1, 2)$ -distributed random variables. Let $U = X - 2Y + 3Z$ and $V = 2X - Y + Z$. Give the joint distribution of U and V .



3.26. For many of the above problems it is instructive to simulate the corresponding model on a computer in order to better understand the theory.

- (a) Generate 10^5 points (X, Y) from the model in Problem 3.6.
- (b) Compare the fraction of points falling in the unit square $[0, 1] \times [0, 1]$ with the theoretical probability in Problem 3.6(b).
- (c) Do the same for the probability $\mathbb{P}(X + Y < 1)$.



3.27. Simulate 10^5 draws from $U(-\pi/2, \pi/2)$ and transform these using the tangent function, as in Problem 3.9. Compare the histogram of the transformed values with the theoretical (Cauchy) pdf.



3.28. Simulate 10^5 independent draws of (U, V) in Problem 3.17. Verify with a histogram of the U -values that the pdf of U is of the form (3.38).



3.29. Consider the MATLAB experiments in Example 3.14.

- (a) Carry out the experiments with $\varrho = 0.4, 0.7, 0.9, 0.99$, and -0.8 , and observe how the outcomes change.
- (b) Plot the corresponding pdfs, as in Fig. 3.6.
- (c) Give also the contour plots of the pdfs, for $\varrho = 0$ and $\varrho = 0.8$. Observe that the contours are *ellipses*.
- (d) Show that these ellipses are of the form

$$x_1^2 + 2\varrho x_1 x_2 + x_2^2 = \text{constant}.$$

Appendix A

Matlab Primer

MATLAB, a portmanteau of MATrix LABoratory, is an interactive matrix-based program for numerical computation. It is a very easy to use high-level language that requires minimal programming skills. The purpose of this appendix is to introduce the reader to some basic MATLAB functions that are used in the main text. For more detailed information and full documentation, please visit the official documentation site

<http://www.mathworks.com/help/techdoc/>.

In addition, the command `help function_name` gives information about the function `function_name`. Alternatively, select Help -> Product Help in the toolbar in the MATLAB command window.

A.1 Matrices and Matrix Operations

The most fundamental objects in MATLAB are, not surprisingly, matrices. For instance, to create a 1×3 matrix (i.e., row vector) **a**, enter into the MATLAB command window:

```
a = [1 2 3]
```

MATLAB returns

```
a =  
    1    2    3
```

To create a matrix with more than one row, use semicolons to separate the rows. For example, the line

```
A = [1 2 3; 4 5 6; 7 8 9];
```

creates a 3×3 matrix A . It is worth noting that MATLAB is case sensitive for variable names and built-in functions. That means MATLAB treats a and A as different objects. To display the i th element in a vector \mathbf{x} , just type $\mathbf{x}(i)$. For example,

```
a(2)
```

refers to the second element of \mathbf{a} . Similarly, one can access a particular element of A by specifying its row and column number (row first followed by column). For instance,

```
A(2,3)
```

displays the $(2, 3)$ entry of the matrix A . To display multiple elements in the matrix, one can use expressions involving colons. For example,

```
A(1,1:2)
```

displays the first and second elements in the first row, whereas

```
A(:,2)
```

displays all the elements in the second column.

To perform numerical computation, one needs some basic matrix operations. In MATLAB, the following matrix operations, among many others, are available:

+	addition	/	right division
−	subtraction	^	power
*	multiplication	'	transpose
\	left division		

For example,

```
a'
```

returns the transpose of **a**:

```
ans =
```

```
    1
    2
    3
```

whereas

```
a*A
```

gives the product of **a** and *A*:

```
ans =
```

```
    30    36    42
```

Other operations are obvious, except for the matrix divisions \backslash and $/$. If *A* is an invertible square matrix and **a** is a compatible vector, then $\mathbf{x} = A \backslash \mathbf{a}$ is the solution of $A \mathbf{x} = \mathbf{a}$ and $\mathbf{x} = \mathbf{a}/A$ is the solution of $\mathbf{x} A = \mathbf{a}$. In other words, $A \backslash \mathbf{a}$ gives the same result (in principle) as $A^{-1} \mathbf{a}$, though they compute their results in different ways. Specifically, the former solves the linear system $A \mathbf{x} = \mathbf{a}$ for **x** by Gaussian elimination, whereas the latter first computes the inverse A^{-1} , then multiplies it by **a**. As such, the second method is in general slower as computing the inverse of a matrix is time-consuming (and inaccurate).

It is important to note that although addition and subtraction are element-wise operations, the other operations listed above are not—they are matrix operations. For example, A^2 gives the square of the matrix *A*, not a matrix whose entries are the squares of those in *A*. One can make the operations $*$, \backslash , $/$, and $^$ to operate element-wise by preceding them by a full stop. For example,

```
A.^2
```

returns the square of the matrix *A*:

```
ans =
```

```
    30    36    42
    66    81    96
   102   126   150
```

On the other hand,

```
A.^2
```

computes the squares element-wise:

```
ans =  
  
     1     4     9  
    16    25    36  
    49    64    81
```

A.2 Some Useful Built-In Functions

In this section we list some common built-in functions which are used throughout the main text. One can learn more about a specific function, say, `eye`, by typing `help eye` in the command window. Here are some useful matrix-building functions:

<code>eye</code>	create an identity matrix
<code>zeros</code>	create a matrix of zeros
<code>ones</code>	create a matrix of ones
<code>diag</code>	create a diagonal matrix or extract the diagonal from a matrix
<code>rand</code>	generate $U(0, 1)$ random variables
<code>randn</code>	generate $N(0, 1)$ random variables

For example, `eye(n)` creates an $n \times n$ identity matrix, and `ones(m,n)` produces an $m \times n$ matrix of ones. Given the 3×3 matrix A ,

```
diag(A)
```

extracts the diagonal of the matrix A :

```
ans =  
  
     1  
     5  
     9
```

But for the 3×1 vector \mathbf{a} , the same command

```
diag(a)
```

builds a diagonal matrix whose main diagonal is **a**:

```
ans =
     1     0     0
     0     2     0
     0     0     3
```

Some other useful vector and matrix functions:

exp	exponential	log	natural log
sqrt	square root	abs	absolute value
sin	sine	cos	cosine
sum	sum	prod	product
max	maximum	min	minimum
chol	Cholesky factorization	inv	inverse
det	determinant	size	size

If **x** is a vector, `sum(x)` returns the sum of the elements in **x**. For a matrix **X**, `sum(X)` returns a row vector consisting of sums of each column, while `sum(X, 2)` returns a column vector of sums of each row. For example,

```
sum(A)
```

returns

```
ans =
    12    15    18
```

whereas

```
sum(A, 2)
```

gives

```
ans =
     6
    15
    24
```

For a positive definite matrix **C**, `chol(C, 'lower')` returns the lower Cholesky factorization **B** such that $BB^T = C$. For example,


```
B = [ 1 0 0; 2 3 0; 4 5 6];
C = B*B';
chol(C,'lower')
```

returns the lower Cholesky factor of BB^T , which is, of course, B .

A.3 Flow Control

MATLAB has the usual control flow statements such as `if-then-else`, `while`, and `for`. For instance, the general form of a simple `if` statement is

```
if condition
    statements
end
```

The statements will be executed if the condition is true. Multiple branching is done by using `elseif` and `else`. For example, the following code simulates rolling a four-sided die:

```
u = rand;
if u <= .25
    disp('1');
elseif u <= .5
    disp('2');
elseif u <= .75
    disp('3');
else
    disp('4');
end
```

The general form of a `while` loop is

```
while condition
    statements
end
```

The statements will be repeatedly executed while the condition remains true. To illustrate the `while` loop syntax, suppose we wish to generate a positive normal random variable [with pdf given in (2.25)]. We can do that using the following simple `while` loop:

```
u = randn;
while u <= 0
    u = randn;
end
```

Another useful control flow statement is the `for` loop, whose general form is

```
for count
    statements
end
```

Unlike a `while` loop, the `for` loop executes the statements for a fixed number of times. As an example, the following code generates five draws from the positive normal distribution.

```
x = zeros(1,5); %% create a storage vector
for i=1:5
    u = randn;
    while u <= 0
        u = randn;
    end
    x(i)=u;
end
```

A.4 Function Handles and Function Files

In previous sections we have introduced some built-in functions in MATLAB. For instance, `sqrt` is a function that takes an argument and returns an output (its square root). Later on we will need to create our own functions that take one or more input arguments, operate on them, and return the results. One way to create new functions is through function handles. For example,

```
f = @(x) x.^2 + 5*x - 10 ; % Note the use of the dot
```

creates the function $f(x) = x^2 + 5x - 10$ with the function handle `f`. The function handle gives you a means of invoking the function. To evaluate, say, $f(10)$, we can type `feval(f,10)`, or simply, `f(10)`.

Function handles can be passed to other functions as inputs. For instance, if we want to find the minimum point of $f(x)$ in the interval $(-10, 10)$, we can use the built-in function `fminbnd` (see Sect. A.6 for a more detailed discussion on optimization routines):

```
[xmin fmin] = fminbnd(f,-10, 10)
```

Note that `fminbnd` takes three inputs (a function handle and the two end-points of the interval) and returns two outputs (the minimizer and the minimal value). In our example, the minimizer of $f(x)$ in $(-10, 10)$ is -2.5 , and the corresponding

functional value is -16.25 . To create a function that takes more than one input is just as easy. For example,

```
g = @(x,y) x.^2 + y.^3 + x.*y
```

defines the two-variable function $g(x, y) = x^2 + y^3 + xy$.

For more complex functions that involve multiple lines and intermediate variables, we need the command `function`. For example, the following code takes a column vector of data and computes its mean and standard deviation:

```
function [meanx, stdevx] = stat(x)
n = length(x);
meanx = sum(x)/n;
stdevx = sqrt(sum(x.^2)/n - meanx.^2);
```

It is important to note that all the code must be written and saved in a separate m-file. Also, the name of the file should coincide with the name of the function; in this case, the file must be called `stat.m`. After saving the file, it can be used the same way as other built-in functions, for example:

```
[meanx stdx] = stat(randn(100,1))
```

returns

```
meanx =
    0.0530

stdx =
    0.9902
```

A.5 Graphics

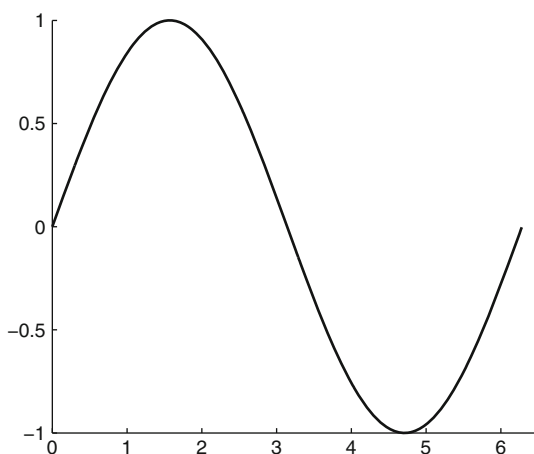
MATLAB has several high-level graphical routines and very extensive plotting capabilities. It allows users to create various graphical objects including two- and three-dimensional graphs. One can also have a title on a graph, add a legend, change

the font and font size, label the axis, etc. For more information, in the command window, click on **Help** and next select **Demos**. Then choose **Graphics** followed by **2D Plots**.

In **MATLAB** the most basic function used to create 2D graphs is `plot`. For example, to make a graph of $y = \sin(x)$ on the interval from $x = 0$ to $x = 2\pi$, we use the following code:

```
x = 0:.01:2*pi;  
y = sin(x);  
plot(x,y);
```

Fig. A.1 A plot of the graph $y = \sin(x)$ from 0 to 2π



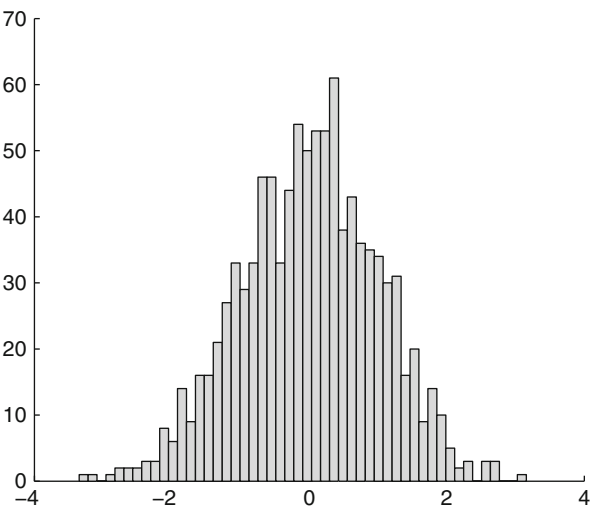
The graph produced is given in Fig. A.1. Note that the command `x = 0:.01:2*pi;` creates a vector whose components range from 0 to 2π in steps of 0.01. Another useful command to create a grid is `linspace` (use `help linspace` to learn more about this function).

Another useful function is `hist`, which allows us to plot histograms. For example,

```
hist(randn(1000,1),50);
```

creates a histogram where the 1000 standard normal draws are put into 50 equally spaced bins (see Fig. A.2).

Fig. A.2 A histogram of 1000 standard normal draws



Instead of a histogram, it is often more useful to have a density estimate. One fast and reliable Gaussian kernel density estimator is the **theta KDE** of Botev et al. (2010). The MATLAB function `kde.m` can be downloaded from <http://www.mathworks.com/matlabcentral/fileexchange/14034-kernel-density-estimator>. See also Example 7.4 for an illustration.

202

It is often desirable to plot several graphs in the same figure window. For this purpose we need the function `subplot(i,j,k)`. The function `subplot(i,j,k)` takes three arguments: the first two tells MATLAB that an $i \times j$ array of plots will be created, and the third is the running index that indicates the k th subplot is currently generated. Suppose we wish to plot the functions $y = \sin(x^2/2)$ and $y = \sin(2x)$ in the same figure window. A little modification of the above code accomplishes this goal:

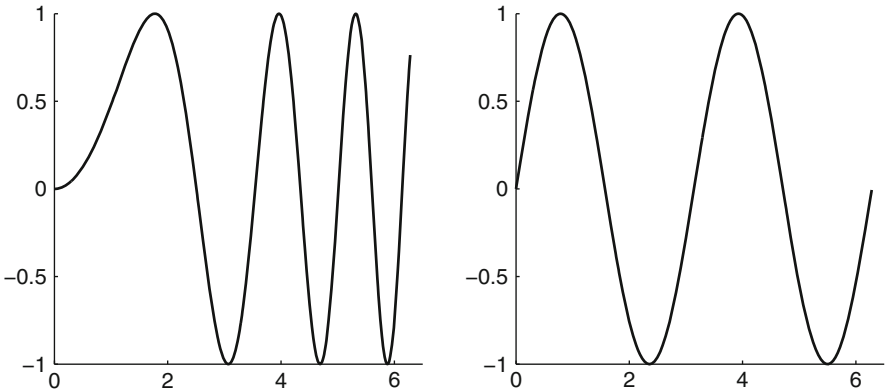


Fig. A.3 Plots of the graphs $y = \sin(x^2/2)$ and $y = \sin(2x)$ from 0 to 2π

```
x = 0:.01:2*pi;
y1 = sin(x.^2/2);    y2 = sin(2*x);
subplot(1,2,1); plot(x,y1);
subplot(1,2,2); plot(x,y2);
```

In addition, one can also easily produce 3D graphical objects in MATLAB. To illustrate various useful routines, suppose we want to plot the density function of the bivariate normal distribution (see Sect. 3.6) given by

 82

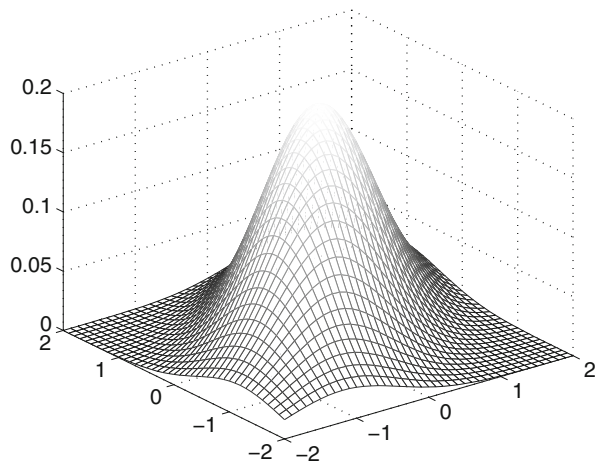
$$f(x, y; \varrho) = \frac{1}{2\pi\sqrt{1-\varrho^2}} e^{-\frac{1}{2(1-\varrho^2)}(x^2-2\varrho xy+y^2)}.$$

As in plotting a 2D graph, we first need to build a grid, and this can be done with the function `meshgrid`. After computing the values of the function at each point on the grid, we can plot the 3D graph using `mesh`. For example, we use the following code

```
rho = .6;
[x y] = meshgrid(-2:.1:2, -2:.1:2); %% build a 2D grid
z = 1/(2*pi*sqrt(1-rho^2)) ...
    * exp(-(x.^2 - 2*rho*x.*y + y.^2)/(2*(1-rho^2)));
mesh(x,y,z);
```

to plot the bivariate normal density function with $\varrho = 0.6$ in Fig. A.4. The ellipsis (...) is used to break up a long line into multiple lines.

Fig. A.4 The density function of the bivariate normal distribution with $\varrho = 0.6$

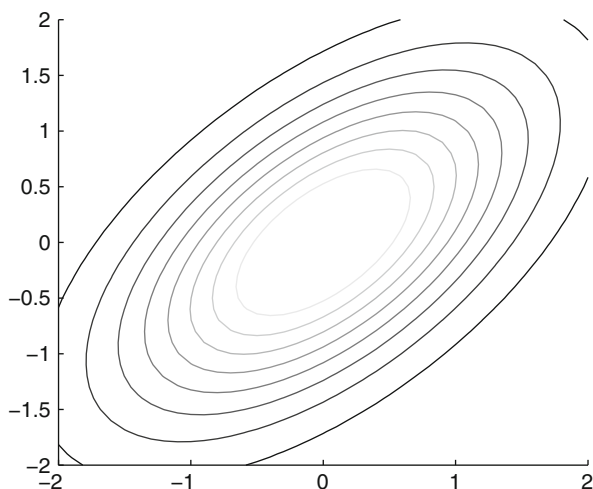


We can produce a contour plot by using the function `contour`:

```
contour(x,y,z);
```

The result is shown in Fig. A.5.

Fig. A.5 A contour plot of the bivariate normal density function with $\rho = 0.6$



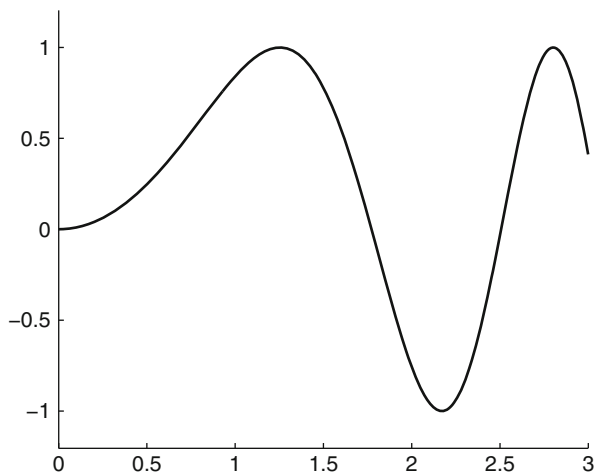
A.6 Optimization Routines

MATLAB provides various built-in optimization routines. In this section we discuss some of them that are used in the main text. Note that all the optimization routines in MATLAB are framed in terms of minimization. In order to perform maximization, some minor changes to the objective function are required. More precisely, suppose we want to maximize the function $f(\mathbf{x})$ and find a maximizer $\mathbf{x}_{\max} = \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$. Instead of the original maximization problem, consider minimizing $-f(\mathbf{x})$ and noting that

$$\mathbf{x}_{\max} = \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x}} -f(\mathbf{x}) .$$

Hence, without loss of generality, we will focus on minimization routines. One basic minimization function is `fminbnd`, which finds the minimum of a single-variable function on a fixed interval. To illustrate its usage, suppose we wish to minimize the function $f(x) = \sin(x^2)$ over the interval $[0, 3]$ (see Fig. A.6).

Fig. A.6 A plot of $f(x) = \sin(x^2)$ from 0 to 3



After defining the function $f(x) = \sin(x^2)$ using the command

```
f = @(x) sin(x.^2);
```

we pass `f` to `fminbnd`, which takes three inputs (the function handle, lower and upper bounds of the interval) and gives two outputs (the minimizer and value of the function evaluated at the minimizer):

```
[xmin fmin] = fminbnd(f,0,3);
```

For this example, we have

```
[xmin fmin]
ans =
    2.1708    -1.0000
```

The function `fminbnd` can only be used to minimize univariate functions on a closed interval. For multivariate minimization, one very useful function is `fminsearch` that finds the unconstrained minimum of a function of several variables. `fminsearch` takes two inputs, namely, the function handle and a starting value. Like `fminbnd`, `fminsearch` gives two outputs: the minimizer and the minimum (the value of the function evaluated at the minimizer). As an example, suppose we wish to maximize the bivariate normal pdf

$$f(x_1, x_2; \varrho) = \frac{1}{2\pi\sqrt{1-\varrho^2}} e^{-\frac{1}{2(1-\varrho^2)}(x_1^2 - 2\varrho x_1 x_2 + x_2^2)}$$

with respect to $\mathbf{x} = (x_1, x_2)$ with $\varrho = 0.6$. To this end, first define $g(x_1, x_2) = -f(x_1, x_2; \varrho)$:

```
rho = .6;
g = @(x) -1/(2*pi*sqrt(1-rho^2)) ...
    *exp(-(x(1).^2 -2*rho*x(1).*x(2) +x(2).^2)
    /(2*(1-rho^2)));
```

Note that the variable \mathbf{x} is a 1×2 vector. Then, we pass g to `fminsearch` with starting values, say, $[1, -1]$:

```
[xmin gmin] = fminsearch(g, [1 -1]);
```

For this example, we have

```
[xmin gmin]
ans =
    0.0000    0.0000   -0.1989
```

That is, the mode of $f(x_1, x_2; \varrho = 0.6)$ is $\mathbf{x} = (0, 0)$, and $f(0, 0) = 0.1989$.

A.7 Handling Sparse Matrices

A **sparse matrix** is simply a matrix that contains a large proportion of zeros. Computation for sparse matrices can typically be done much faster than for full matrices. In addition, as most of the elements in a sparse matrix are zeros, the storage cost of a sparse matrix is also small. In statistics we often need to deal with large sparse matrices. Thus it is useful to learn how to handle them in MATLAB.

A basic function for creating sparse matrices is `sparse`. For example, suppose the matrix

$$W = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 3 & 1 \end{pmatrix}$$

is stored as a full matrix in MATLAB. The command `sparse(W)` converts W to sparse form by squeezing out any zero elements and returns:

```
ans =  
  
    (1,1)      1  
    (2,2)      1  
    (3,3)      2  
    (4,4)      3  
    (4,5)      1
```

Notice that only the nonzero elements in W are stored. In general, we can create a matrix S by the command `S = sparse(i,j,s,m,n)`, which uses vectors \mathbf{i} , \mathbf{j} , and \mathbf{s} to generate an $m \times n$ sparse matrix such that $S(\mathbf{i}(k), \mathbf{j}(k)) = \mathbf{s}(k)$. For example, to create the matrix W above, we first need to build a vector \mathbf{s} that stores all the nonzero elements:

```
s = [1 1 2 3 1]';
```

Next, we create a vector \mathbf{i} that stores the row position for each element in \mathbf{s} . For example, the first element in \mathbf{s} should be in the first row, the second element in second row, and so on. We then do the same thing for the column positions and store them in the vector \mathbf{j} :

```
i = [1 2 3 4 4]';  
j = [1 2 3 4 5]';
```

Finally,

```
W = sparse(i,j,s,4,5);
```

creates the 4×5 matrix W above.

There are several useful built-in functions for creating special sparse matrices. For example,

```
I = speye(100);
```

creates the 100×100 sparse identity matrix. Of course we can accomplish the same goal by using

```
I = sparse(1:100,1:100,ones(1,100));
```

though the latter is more clumsy. Another useful function is `spdiags`, the sparse version of `diag`, which can be used to extract and create sparse diagonal matrices. Use `help spdiags` to learn more about this function.

As mentioned earlier, one main advantage of working with sparse rather than full matrices is that computations involving sparse matrices are usually much quicker. For instance, it takes about 2.7 seconds to obtain the Cholesky decomposition of the full 5000×5000 identity matrix:

```
tic; chol(eye(5000)); toc;
Elapsed time is 2.728245 seconds.
```

whereas the same operation takes only 0.015 second for a sparse 5000×5000 identity matrix:

```
tic; chol(speye(5000)); toc;
Elapsed time is 0.014867 seconds.
```

A.8 Gamma and Dirichlet Generator

The following MATLAB program `gamrand` implements the method developed in (Marsaglia and Tsang 2000) to generate samples from a $\text{Gamma}(\alpha, \lambda)$ distribution. If the *Statistics Toolbox* is available, the function `gamrnd` can be used instead; but note that `gamrnd(a,b)` generates random variables from a $\text{Gamma}(a, 1/b)$ distribution.

```
function x=gamrand(alpha,lambda)
if alpha>1
    d=alpha-1/3; c=1/sqrt(9*d); flag=1;
    while flag
        Z=randn;
        if Z>-1/c
            V=(1+c*Z)^3; U=rand;
            flag=log(U)>(0.5*Z^2+d-d*V+d*log(V));
        end
    end
end
```

```

        x=d*V/lambda;
else
    x=gamrand(alpha+1,lambda);
    x=x*rand^(1/alpha);
end

```

As a direct consequence of Theorem 8.2, the following MATLAB program `dirichrnd` generates samples from a Dirichlet(α) distribution. Draws from a Beta(α, β) are obtained by taking $\alpha = (\alpha, \beta)$.

 241

```

function x=dirichrnd(alpha)
n=length(alpha)-1;
Y=nan(1,n+1);
for k=1:n+1
    Y(k)=gamrand(alpha(k),1);
end
x=Y(1:n)/sum(Y);

```

A.9 Cdfs and Inverse Cdfs

The following MATLAB program `cumcdf` evaluates the cdfs of normal, Student's t , gamma, chi-squared, and F distributions. If the *Statistics Toolbox* is available, the function `cdf` can be used instead.

```

function y = cumcdf(dist,x,varargin)
switch dist
case 'norm'
    mu = varargin{1}; sigma = varargin{2};
    y = (erf((x - mu)/sigma)/sqrt(2)) + 1)/2;
case 't'
    nu = varargin{1};
    y = 1-0.5*betainc(nu/(nu+x.^2),nu/2,1/2);
case 'gamma'
    alpha = varargin{1}; lambda = varargin{2};
    % different from Stats toolbox
    y = gammainc(lambda*x,alpha);
case 'chi2'
    n = varargin{1};
    y = gammainc(x/2,n/2);
case 'F'
    m = varargin{1}; n = varargin{2};
    y = 1 - betainc(n/(n+m*x),n/2,m/2);
end

```

The following MATLAB program `icumdf` evaluates the inverse cdfs of normal, Student's t , gamma, chi-squared, and F distributions. The corresponding built-in function in the *Statistics Toolbox* is `icdf`.

```
function x = icumdf(dist,y,varargin)
switch dist
    case 'norm'
        mu = varargin{1}; sigma = varargin{2};
        x = mu + sigma*sqrt(2)*erfinv(2*y -1);
    case 't'
        nu = varargin{1};
        x = sqrt(nu/betaincinv(2*(1-y),nu/2,1/2) - nu);
    case 'gamma'
        alpha = varargin{1}; lambda = varargin{2};
        % different from Stats toolbox
        x = gammaincinv(y,alpha)/lambda;
    case 'chi2'
        n = varargin{1};
        x = gammaincinv(y,n/2)*2;
    case 'F'
        m = varargin{1}; n = varargin{2};
        x = n/m/betaincinv(1-y,n/2,m/2) - n/m;
end
```

A.10 Further Reading and References

The official MATLAB documentation site is

<http://www.mathworks.com/help/techdoc/>

A good place to learn more about the major functionality in MATLAB is the MATLAB *Getting Started Guide* available at

http://www.mathworks.com/help/pdf_doc/matlab/getstart.pdf

MATLAB programs for generating random variables from a wide range of distributions may be found on the homepage of the *Handbook of Monte Carlo Methods* (Kroese et al. 2011):

<http://www.montecarlohandbook.org>

Finally, all programs and (large) data files in this book may be downloaded from the homepage

<http://www.statmodcomp.org>

To accommodate the users of the statistical programming language R, we have mirrored each MATLAB program with its equivalent in R.

Appendix B

Mathematical Supplement

B.1 Multivariate Differentiation

For a real-valued multivariate function $f(x_1, \dots, x_n)$ the **partial derivative** with respect to x_i , denoted $\frac{\partial f}{\partial x_i}$ or simply $\partial_i f$, is the derivative taken with respect to x_i while all other variables are held constant. The partial derivative of $\partial_i f$ with respect to x_j is denoted $\frac{\partial^2 f}{\partial x_i \partial x_j}$ or simply $\partial_{ij} f$.

Let \mathbf{f} be a multivariate function taking values in \mathbb{R}^m , defined by

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \mapsto \begin{pmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{pmatrix} = \mathbf{f}(\mathbf{x}) .$$

The **derivative** of \mathbf{f} at \mathbf{x} is defined as the matrix of partial derivatives,

$$J_{\mathbf{f}}(\mathbf{x}) = \begin{pmatrix} \partial_1 f_1(\mathbf{x}) & \cdots & \partial_n f_1(\mathbf{x}) \\ \vdots & \cdots & \vdots \\ \partial_1 f_m(\mathbf{x}) & \cdots & \partial_n f_m(\mathbf{x}) \end{pmatrix} , \quad (\text{B.1})$$

and is called the **matrix of Jacobi** of \mathbf{f} at \mathbf{x} , sometimes written as $\frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x})$.

Example B.1 (Differentiating a Linear Function). Let $\mathbf{f}(\mathbf{x}) = A\mathbf{x}$ for some $m \times n$ constant matrix A . Then,

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = A . \quad (\text{B.2})$$

To see this, let a_{ij} denote the (i, j) th element of A , so that

$$\mathbf{f}(\mathbf{x}) = A\mathbf{x} = \begin{pmatrix} \sum_{k=1}^n a_{1k}x_k \\ \vdots \\ \sum_{k=1}^n a_{mk}x_k \end{pmatrix}.$$

To find the (i, j) th element of the $m \times n$ Jacobian matrix \mathbf{J}_f , we differentiate the i th element of \mathbf{f} with respect to x_j :

$$\frac{\partial f_i(\mathbf{x})}{\partial x_j} = \frac{\partial}{\partial x_j} \sum_{k=1}^n a_{ik}x_k = a_{ij}.$$

In other words, the (i, j) th element of \mathbf{J}_f is a_{ij} , the (i, j) th element of A .

For a real-valued multivariate function, that is, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the **gradient** of f is the transpose of the Jacobian matrix, that is, the *column* vector

$$\nabla f(\mathbf{x}) = \begin{pmatrix} \partial_1 f(\mathbf{x}) \\ \vdots \\ \partial_n f(\mathbf{x}) \end{pmatrix}. \quad (\text{B.3})$$

The derivative of the function $\mathbf{x} \mapsto \nabla f(\mathbf{x})$ is called the **Hessian matrix** of f , denoted $H_f(\mathbf{x})$ or $\nabla^2 f(\mathbf{x})$. In other words, the Hessian is the matrix of second derivatives:

$$\nabla^2 f(\mathbf{x}) = \begin{pmatrix} \partial_{11} f(\mathbf{x}) & \cdots & \partial_{1n} f(\mathbf{x}) \\ \vdots & \cdots & \vdots \\ \partial_{n1} f(\mathbf{x}) & \cdots & \partial_{nn} f(\mathbf{x}) \end{pmatrix}. \quad (\text{B.4})$$

If the partial derivatives are *continuous* in a region around \mathbf{x} , then $\partial_{ij} f(\mathbf{x}) = \partial_{ji} f(\mathbf{x})$ and, hence, the Hessian matrix $H_f(\mathbf{x})$ is *symmetric*.

Example B.2 (Differentiating a Quadratic Function). Let $f(\mathbf{x}) = \mathbf{x}^\top A\mathbf{x}$ for some $n \times n$ constant matrix A . Then,

$$\nabla f(\mathbf{x}) = (A + A^\top)\mathbf{x}. \quad (\text{B.5})$$

It follows immediately that if A is *symmetric*, i.e., $A = A^\top$, then $\nabla(\mathbf{x}^\top A\mathbf{x}) = 2A\mathbf{x}$ and $\nabla^2(\mathbf{x}^\top A\mathbf{x}) = 2A$.

To prove (B.5), first, observe that the quadratic function $f(\mathbf{x}) = \mathbf{x}^\top A\mathbf{x}$ is real-valued, and therefore the Jacobian \mathbf{J}_f is a $1 \times n$ vector (and its transpose is the gradient). Specifically,

$$f(\mathbf{x}) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j ,$$

and the k th element of \mathbf{J}_f is obtained by differentiating $f(\mathbf{x})$ with respect to x_k :

$$\frac{\partial f(\mathbf{x})}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j = \sum_{j=1}^n a_{kj} x_j + \sum_{i=1}^n a_{ik} x_i .$$

The first term on the right-hand side is equal to the k th element of $\mathbf{A}\mathbf{x}$, whereas the second term equals the k th element of $\mathbf{x}^\top \mathbf{A}$, or equivalently the k th element of $\mathbf{A}^\top \mathbf{x}$.

Gradients and Hessian matrices feature prominently in multidimensional Taylor expansions.

Theorem B.1. (Multidimensional Taylor Expansions). Let \mathcal{X} be an open subset of \mathbb{R}^n and let $\mathbf{a} \in \mathcal{X}$. If $f : \mathcal{X} \rightarrow \mathbb{R}$ is a continuously twice differentiable function with gradient $\nabla f(\mathbf{x})$ and Hessian matrix $H_f(\mathbf{x})$, then for every $\mathbf{x} \in \mathcal{X}$ we have the following first- and second-order Taylor expansions:

$$f(\mathbf{x}) = f(\mathbf{a}) + [\nabla f(\mathbf{a})]^\top (\mathbf{x} - \mathbf{a}) + \mathcal{O}(\|\mathbf{x} - \mathbf{a}\|^2)$$

and

$$f(\mathbf{x}) = f(\mathbf{a}) + [\nabla f(\mathbf{a})]^\top (\mathbf{x} - \mathbf{a}) + \frac{1}{2} (\mathbf{x} - \mathbf{a})^\top H_f(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \mathcal{O}(\|\mathbf{x} - \mathbf{a}\|^3)$$

as $\|\mathbf{x} - \mathbf{a}\| \rightarrow 0$. By dropping the \mathcal{O} remainder terms, one obtains the corresponding Taylor approximations.

B.2 Proof of Theorem 2.6 and Corollary 2.2

 34

The proof makes use of two fundamental properties of the expectation \mathbb{E} : the *monotone convergence theorem* and the *dominated convergence theorem*. The first states that if $X_1 \leq X_2 \leq X_3 \leq \dots$ is a sequence of positive random variables that increases to a random variable X , then the corresponding expectations $\mathbb{E}X_1 \leq \mathbb{E}X_2 \leq \mathbb{E}X_3 \dots$ converge to $\mathbb{E}X$. The second theorem states that the same holds true for any positive sequence X_1, X_2, \dots converging to X , if there exists a Y with $\mathbb{E}Y < \infty$ such that $X_n \leq Y$ for all n . An accessible account of these theorems may be found, for example, in (Williams 1991).

We prove Theorem 2.6 for the case $k = 1$ only. Let $G(z) = \mathbb{E}z^X$. Take a fixed z with $|z| < R$ and any $r < R$ such that $r < |z| < R$. Let (h_n) be any sequence converging to 0 such that $|z + h_n| < r$. By definition, the derivative of G at z is $\lim_{n \rightarrow \infty} \mathbb{E}C_n$, where $C_n = h_n^{-1}[(z + h_n)^X - z^X]$. Observe that

1. $|C_n|$ is dominated by $X r^{X-1}$,
2. $\mathbb{E}X r^{X-1} < \infty$, because the power series $\sum_{x=0}^{\infty} x z^{x-1} f(x)$ has again radius of convergence R ,
3. $\lim_{n \rightarrow \infty} C_n = X z^{X-1}$.

It follows by the dominated convergence theorem that

$$\lim_{n \rightarrow \infty} \mathbb{E}C_n = \mathbb{E} \lim_{n \rightarrow \infty} C_n = \mathbb{E}X z^{X-1}.$$

Next, let (z_n) be a sequence of real numbers that is converging to 1, where $|z_n| < 1$ for all n . The sequence of random variables (Y_n) defined by $Y_n = X(X-1)\cdots(X-k+1)z_n^k$ is increasing to $Y = X(X-1)\cdots(X-k+1)$. Hence, by the monotone convergence theorem $\lim_{n \rightarrow \infty} \mathbb{E}Y_n = \mathbb{E}Y$. This shows (2.10). The second statement of the corollary is left as an exercise.

B.3 Proof of Theorem 2.7

If the moment generating function of a random variable X is finite in an open interval containing 0, then for all $n = 0, 1, \dots$,

$$\mathbb{E}X^n = M^{(n)}(0),$$

where $M^{(n)}$ is the n th derivative of the MGF M evaluated at 0.

Proof. Let $R > 0$ be such that $M(s) < \infty$ for all $|s| < R$. Choose any numbers r and s such that $0 < r < R$ and $|s| < r$. Let (h_n) be a sequence converging to 0 satisfying $|h_n| < \varepsilon$ and $|s + h_n| < r$ for some $\varepsilon > 0$. Let $C_n = h_n^{-1}[e^{(s+h_n)X} - e^{sX}] = e^{sX}(e^{h_nX} - 1)/h_n$, which converges to Xe^{sX} . Also, $|C_n| \leq H(X) \stackrel{\text{def}}{=} e^{(|s|+\varepsilon)|X|}|X|$, because $0 \leq (e^t - 1)/t \leq e^{|t|}$ for all t . Moreover, because $|s| + \varepsilon < r$ and x grows at a lesser rate than e^{ax} for any $a > 0$, there must exist an $M > 0$ such that for all $|x| > M$, $H(x) < e^{r|x|}$. It follows that

$$\begin{aligned} \mathbb{E}H(X) &\leq \mathbb{E}H(X) \mathbf{I}_{\{|X| > M\}} + \mathbb{E}H(X) \mathbf{I}_{\{|X| \leq M\}} \\ &\leq \mathbb{E}e^{r|X|} + \max_{|x| \leq M} H(x) < \infty. \end{aligned}$$

By the dominated convergence theorem we have $M'(s) = \lim_{n \rightarrow \infty} \mathbb{E}C_n = \mathbb{E} \lim_{n \rightarrow \infty} C_n = \mathbb{E}[Xe^{sX}]$. Finally, take a monotone sequence (s_n) converging to 0 and apply the monotone convergence theorem to the sequence (Xe^{s_nX}) to find $M'(0) = \mathbb{E}X$. The proof for higher moments is similar. \square

B.4 Proof of Theorem 3.10

Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ be an orthonormal basis of \mathbb{R}^n such that $\mathbf{v}_1, \dots, \mathbf{v}_k$ spans \mathcal{V}_k and $\mathbf{v}_1, \dots, \mathbf{v}_m$ spans \mathcal{V}_m . We can write the orthogonal projection matrices onto \mathcal{V}_j , as $P_j = \sum_{i=1}^j \mathbf{v}_i \mathbf{v}_i^\top$, $j = k, m, n$, where \mathcal{V}_n is defined as \mathbb{R}^n . Note that P_n is simply the identity matrix. Let $V = (\mathbf{v}_1, \dots, \mathbf{v}_n)$ and define $\mathbf{Z} = (Z_1, \dots, Z_n)^\top = V^\top \mathbf{X}$. Recall that any orthogonal transformation such as $\mathbf{z} = V^\top \mathbf{x}$ is *length preserving*; that is, $\|\mathbf{z}\| = \|\mathbf{x}\|$.

To prove the first statement of the theorem, note that $V^\top \mathbf{X}_j = V^\top P_j \mathbf{X} = (Z_1, \dots, Z_j, 0, \dots, 0)^\top$, $j = k, m$. It follows that $V^\top (\mathbf{X}_m - \mathbf{X}_k) = (0, \dots, 0, Z_{k+1}, \dots, Z_m, 0, \dots, 0)^\top$ and $V^\top (\mathbf{X} - \mathbf{X}_m) = (0, \dots, 0, Z_{m+1}, \dots, Z_n)^\top$. Moreover, being a linear transformation of a normal random vector, \mathbf{Z} is also normal, with covariance matrix $V^\top V = I$; see also Problem 3.13. In particular, the $\{Z_i\}$ are *independent*. This shows that \mathbf{X}_k , $\mathbf{X}_m - \mathbf{X}_k$, and $\mathbf{X} - \mathbf{X}_m$ are independent as well. 95

Next, observe that $\|\mathbf{X}_k\| = \|V^\top \mathbf{X}_k\| = \|\mathbf{Z}_k\|$, where $\mathbf{Z}_k = (Z_1, \dots, Z_k)^\top$. The latter vector has independent components with variances 1, and its squared norm has therefore (by definition) a $\chi_k^2(\theta)$ distribution. The noncentrality parameter is $\theta = \|\mathbb{E} \mathbf{Z}_k\| = \|\mathbb{E} \mathbf{X}_k\| = \|\boldsymbol{\mu}_k\|$, again by the length-preserving property of orthogonal transformations. This shows that $\|\mathbf{X}_k\|^2 \sim \chi_k^2(\|\boldsymbol{\mu}_k\|)$. The distributions of $\|\mathbf{X}_m - \mathbf{X}_k\|^2$ and $\|\mathbf{X} - \mathbf{X}_m\|^2$ follow by analogy. \square

B.5 Proof of Theorem 5.2

First, observe that, by Theorem 5.1, 131

$$\frac{(m-1)S_X^2}{\sigma^2} \sim \chi_{m-1}^2 \quad \text{and} \quad \frac{(n-1)S_Y^2}{\sigma^2} \sim \chi_{n-1}^2.$$

Because these random variables are independent of each other, their sum, V , say, can be written as the sum of $m+n$ independent squared standard normal random variables and has therefore a χ_{m+n-2}^2 distribution. Thus,

$$V = \frac{(m+n-2)S_p^2}{\sigma^2} \sim \chi_{m+n-2}^2.$$

Second, let

$$Z = \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sigma / \sqrt{\frac{1}{m} + \frac{1}{n}}}.$$

Then, $Z \sim N(0, 1)$ and the square of the pivot T in Theorem 5.2 can be written as

$$T^2 = \frac{Z^2}{V/(m+n-2)},$$

where Z and V are independent, because \bar{X} and \bar{Y} are independent of each other, and are both independent of S_X^2 and S_Y^2 ; see Theorem 5.1. The random variable T^2 is thus the independent quotient of a χ_1^2 and a χ_{m+n-2}^2 random variable. Hence, by Theorem 3.11, $T^2 \sim F(1, m+n-2)$. It follows now from Theorem 2.19 (and the fact that the pdf of T is symmetric around 0) that $T \sim t_{m+n-2}$. \square

88

51

Index

Symbols

\sim distributed as, 28
 \mathbb{E} expectation, 29
 iid
 \sim independent and identically distributed as, 66
 I indicator, 73
 \cap intersection, 7
 \mathbb{P} probability, 9
 \propto proportional to, 215
 φ standard normal pdf, 46
 Φ standard normal cdf, 46
 \cup union, 7

A

acceptance–rejection method, 55, 214, 215
affine transformation, 47, 75, 76, 82, 83
Akaike information criterion, 305, 320
alternative hypothesis, 140
Analysis of Variance (ANOVA), 111, 142, 143, 156
 model, 111–114
 single-factor, 112, 115, 143
 two-factor, 113
autocorrelation, 290
autocovariance, 290
autoregressive moving average, 287, 303
auxiliary mixture sampling, 340
auxiliary variable methods, 183

B

bag of words method, 261
balanced design, 112
bandwidth, 201
`bar.m`, 4

Bayes factor, 140, 251
 Savage–Dickey density ratio, 254
Bayes’ rule, 16, 227, 228
Bayesian information criterion, 305, 320
Bayesian network, 244–248
Bayesian statistics, 121, 228, 233
belief net, 246
Bernoulli
 distribution, 36
 process, 66
 regression, 266
beta distribution, 74, 229, 241, 256, 365
beta function, 74
bias, 122, 205
binomial distribution, 18, 24, 37, 67, 69, 92
 normal approximation to, 92
binomial formula, 38
binomial model, 135
 two-sample, 103, 136
birthday problem, 15
blocking, 114
bootstrap method, 128, 203, 205
Box–Muller method, 82
burn-in, 213, 291

C

categorical variable, 111
Cauchy distribution, 50, 71, 94, 163, 204
Cauchy–Schwarz inequality, 95, 171
`ceil.m`, 21
central limit theorem, 90, 130
 for random vectors, 92
chi-squared distribution, 48, 86, 89, 96, 132, 134
classical statistics, 121

- coefficient of determination, 156
 - coin tossing, 3, 7, 18, 24, 37, 39, 66, 121, 228
 - combined multiple-recursive generator, 52
 - complete-data likelihood, 183
 - completing the squares, 238, 393
 - concentration matrix, 308
 - conditional
 - expectation, 77
 - pdf, 71
 - probability, 12–19
 - confidence
 - set, 174
 - confidence interval, 128, 174, 175, 206
 - approximate, 128
 - approximate – for p (binomial distribution), 136
 - approximate – for p (two-sample, binomial distribution), 136
 - Bayesian, 128, 229
 - bootstrap, 206
 - for $\mu_X - \mu_Y$ (two-sample normal distribution), 134, 157
 - for σ^2 (normal distribution), 132
 - for σ_X^2/σ_Y^2 (two-sample normal distribution), 134
 - conjugate family, 249
 - consistent estimator, 176
 - convex function, 33
 - correlation coefficient, 76, 85, 95, 124
 - sample, 125, 156
 - counting problems, 19
 - covariance, 75
 - matrix, 77, 79, 83, 84, 86, 92, 168, 285, 306, 308, 310
 - method, 290
 - covariate, 105
 - coverage probability, 128
 - Cramér–Rao inequality, 171
 - credible interval, 128, 229
 - cross-validation, 146
 - K -fold, 147
 - leave-one-out, 147
 - linear model, 148
 - cumdf.m, 60, 366
 - cumsum.m, 4, 55
 - cumulative distribution function (cdf), 25, 29
 - joint, 63
- D**
- data
 - reduction, 150
 - transformation, 110
 - data augmentation, 278
 - De Morgan’s rules, 8, 19
 - delta method, 92, 207
 - dependent variable, 105
 - derivatives
 - multidimensional, 367
 - partial, 367
 - design matrix, 115, 116, 125, 127, 148, 173, 237, 265, 291, 304, 316
 - detailed balance equations, 213, 214
 - digamma function, 191
 - directed acyclic graph, 244
 - Dirichlet distribution, 241, 365
 - discrete joint pdf, 64
 - discrete random variable, 111
 - disjoint events, 7, 9
 - distribution
 - Bernoulli, 36
 - beta, 74, 229, 241, 256, 365
 - binomial, 37, 67, 69, 92
 - Cauchy, 50, 71, 94, 163, 204
 - chi-squared, 48, 86, 89, 96, 132, 134
 - continuous joint, 69, 73
 - Dirichlet, 241, 365
 - discrete joint, 64–69
 - discrete uniform, 59
 - double exponential, 190
 - exponential, 43, 94
 - exponential family, 152, 167, 174, 266
 - F , 50, 51, 89, 134
 - gamma, 48, 49, 232, 242
 - Gaussian, *see* normal
 - geometric, 38
 - inverse-gamma, 234, 318, 333, 336, 342
 - logistic, 59
 - mixed joint, 73–74
 - mixture, 187, 201, 221
 - multinomial, 68, 185, 220, 240
 - multivariate normal, 83, 106, 307
 - multivariate Student’s t , 271, 285
 - noncentral χ^2 , 88
 - normal, 45, 57, 81, 82
 - Poisson, 34, 40
 - positive normal, 56, 71, 354–355
 - Student’s t , 50, 89, 131, 133
 - truncated normal, 279, 286
 - uniform, 42, 188
 - Weibull, 61, 190, 200
 - dominated convergence theorem, 370
 - double exponential distribution, 190
 - drawing with or without replacement, 19
- E**
- efficient score, 167
 - erf.m, 60
 - EM-algorithm, 182, 279, 327

empirical cdf, [196](#), [203](#)
 reduced, [199](#)
 ergodic Markov chain, [212](#)
 error terms, [115](#), [173](#)
 estimate, [122](#)
 estimator, [122](#)
 bias, [122](#)
 unbiased, [122](#)
 event, [6](#)
 elementary, [10](#)
 expectation, [31](#), [29–33](#)
 conditional, [77](#)
 for joint distributions, [74](#)
 properties, [33](#), [75](#)
 vector, [77](#), [79](#), [83](#)
 explanatory variable, [105](#)
 exponential distribution, [43](#), [94](#)
 exponential family, [152](#), [167](#), [174](#), [266](#)
 conjugate prior, [249–251](#)
 information matrix, [170](#)
 natural, [152](#)
 exponential model, [109](#)

F
 factor level, [111](#)
 factorial experiment, [111](#)
 factorization theorem, [150](#)
 F distribution, [50](#), [51](#), [89](#), [134](#)
`find.m`, [55](#)
 Fisher information matrix, [168](#)
 observed, [268](#)
 Fisher's scoring method, [180](#), [283](#)
 full rank matrix, [126](#)
 functions of random variables, [78](#)
`fzero.m`, [193](#)

G
 Galton, Francis, [104](#)
 gamma distribution, [48](#), [49](#), [232](#), [242](#)
 gamma function, [48](#), [49](#), [74](#), [191](#), [193](#)
`gamrand.m`, [232](#), [364](#)
`gamrnd.m`, [337](#), [364](#)
 Gaussian distribution, *see* normal distribution
 generalized likelihood ratio, [178](#)
 generalized linear model, [265](#)
 geometric distribution, [18](#), [38](#)
 geometric sum, [39](#)
 Gibbs sampler, [218–219](#), [225](#), [226](#), [230](#), [232](#),
 [234–236](#), [258–259](#), [280](#), [316–320](#),
 [332–333](#), [335–339](#), [342–345](#)
 global balance equations, [212](#)
 goodness of fit test, [220](#)

gradient, [368](#)
 grid search, [193](#)

H

Hessian matrix, [170](#), [176](#), [180](#), [183](#), [368](#)
 hierarchical model, [229](#), [332](#)
 hyperparameter, [245](#)
 hypothesis testing, [140–195](#)

I

`icumd.f.m`, [60](#), [129](#), [366](#)
 improper prior, [236](#)
 independence
 of events, [17](#)
 of random variables, [65](#), [66](#), [70](#), [75](#)
 independence sampler, [215](#)
 independent and identically distributed (iid),
 [66](#), [71](#), [89](#), [101–104](#), [130](#)
 independent variable, [105](#)
 indicator, [57](#), [74](#)
 initial distribution, [210](#)
 integrated moving average, [301](#)
 interval estimate, *see* confidence interval, [174](#)
 inverse-gamma distribution, [234](#), [318](#), [333](#),
 [336](#), [342](#)
 inverse-transform method, [54](#), [71](#), [200](#), [203](#)
 discrete, [54](#)
 irreducible, [213](#)

J

Jacobian matrix, *see* matrix of Jacobi
 Jensen's inequality, [33](#), [192](#)
 joint
 cdf, [63](#)
 distribution, [63](#), [79](#)
 joint pdf, [69](#)
 for dependent random variables, [67](#)
 jointly normal distribution, *see* multivariate
 normal distribution

K

Kalman filter, [325](#)
`kde.m`, [202](#)
 kernel density estimation, [201–203](#), [209](#), [216](#),
 [232](#)
 Kolmogorov–Smirnov statistic, [200](#), [221](#)
 Kronecker product, [116](#), [314](#), [317](#), [389](#)
 Kullback–Leibler distance, [192](#)

L

Langevin Metropolis–Hastings sampler, 224
 latent variable methods, *see* auxiliary variable methods
 law of large numbers, 89, 130
 law of total probability, 16
 least-squares method, 125–128, 222
 likelihood, 123, 161

- Bayesian, 228
- binomial, 161
- complete-data, 183
- concentrated, 295
- normal, 162
- optimization, 182
- profile, 189, 295, 304

 limiting pdf, 212
 linear model, 115, 173
 linear regression model, 108
 linear transformation, 79
 local balance equations, *see* detailed balance equations
 location family, 171, 182
 log-likelihood, 165
 logistic distribution, 59, 267
 logistic model, 109
 logistic regression, 267
 logit model, 267

M

marginal likelihood, 252
 marginal pdf, 65, 70, 86, 230, 241, 257
 Markov

- property, 209

 Markov chain, 209–213, 217–219, 259, 323

- ergodic, 212
- reversible, 212

 Markov chain Monte Carlo, 209–220, 274, 276, 291

MATLAB

basic matrix operations, 349–352
 built-in functions, 352–354
 for-loop, 355
 function, 356
 function handle, 355
 graphics, 356–360
 if-then-else, 354
 optimization routines, 360–362
 sparse matrix routines, 362–364
 while-loop, 354–355
 matrix

- covariance, 77, 84, 86, 92, 168, 285, 306, 308, 310

 matrix of Jacobi, 81, 242, 257, 284, 367

maximum likelihood estimator, 172–180, 182
 mean square error, 154, 205
 measurement equation, 323
 median, 222

- sample, 204

 memoryless property, 40, 44, 58
 method of moments, 123, 124
 Metropolis–Hastings algorithm, 214–217
 mixture distribution, 187, 201, 221
 mixture model, 187–188
 mode, 172, 229
 model

- Analysis of Variance (ANOVA), 111–114
- autoregressive moving average, 287, 303
- binomial, 103, 135
- exponential, 109
- hierarchical Bayesian model, 229, 332
- linear regression, 108
- logistic, 109
- multinomial, 240
- multiple linear regression, 106, 115
- nested, 253
- normal linear, 88, 114–117, 125, 137, 142, 148, 156, 237
- power law, 109
- probability, 10, 121
- randomized block design, 143
- regression, 104–111
- response surface, 109
- selection, 114, 142, 146, 251, 287
- simple linear regression, 106, 115, 127, 139
- single-factor ANOVA, 112, 143
- state space, 323
- stochastic volatility, 339–345
- time-varying parameter autoregressive, 333–339
- two-factor ANOVA, 113
- unobserved components, 325–333
- Weibull, 109
- zero inflated Poisson, 258

 moment, 32

- sample-, 123

 moment generating function (MGF), 35, 86, 91, 96
 Monte Carlo

- integration, 130
- sampling, 195–226

 Monty Hall problem, 13
 moving average, 289, 297

- integrated, 301

 multinomial distribution, 68, 185, 220, 240
 multinomial model

- Bayesian, 240

multiple linear regression, 106, 115
 multivariate normal distribution, 83, 82–89, 95,
 106, 307

N

natural exponential family, 152
 neighborhood structure, 224
 nested model, 253
 Newton's binomial formula, 38
 Newton–Raphson method, 180
 noncentral χ^2 distribution, 88
 nonlinear regression, 109, 189, 222
 normal distribution, 45, 57, 81, 83
 generating from, 82
 positive, 56, 71, 354–355
 normal equations, 126
 normal linear model, 88, 114–117, 125, 137,
 142, 148, 156, 266
 Bayesian, 237
 normal model
 two-sample, 103, 111, 133
 nuisance factor, 114
 null hypothesis, 140

O

observed information matrix, 268
 orthogonal matrix, 95

P

p -value, 140, 195
 partial derivative, 367
 partition, 16
 Pearson's height data, 104
 pivot variable, 129
 plot.m, 4
 Poisson distribution, 34, 40
 Poisson regression, 282
 polynomial regression, 108
 pooled sample variance, 133
 positive normal distribution, 56, 71, 354–355
 posterior
 mean, 229
 mode, 229
 posterior pdf, 121
 asymptotic normality, 248
 power law model, 109
 precision matrix, 308
 predicted residual, 147
 predictive pdf, 261
 predictor, 105

prior pdf, 227, 249
 improper, 236
 uninformative, 233
 probability, 3, 4, 9–11
 probability density function (pdf)
 discrete joint, 64
 conditional, 67
 continuous, 28
 discrete, 27
 probability distribution, 25
 continuous, 28
 discrete, 27
 probability generating function (PGF), 34
 probability model, 10, 121
 probit model, 273
 product rule, 14, 67, 72, 210, 229, 245
 profile likelihood, 189, 295, 304
 projection matrix, 96, 126, 148
 pseudo-inverse, 126, 239

Q

quad.m, 60
 quotient of independent random variables, 71

R

radius of convergence, 34
 rand.m, 4, 71
 randn.m, 57
 random
 experiment, 3, 5, 10
 number generator, 52
 vector, 79
 random variable, 23
 continuous, 25, 28
 discrete, 25, 111
 functions of, 78
 quotient of, 71
 range, 25
 random vector, 63
 transformation, 81
 random walk sampler, 216
 randomized block design, 143
 range
 of a random variable, 25
 rank, 126
 ratio estimator, 93, 207
 reduction of data, 150
 regression
 line, 106
 model, 104–111
 multiple linear, 106

regression (*cont.*)
 nonlinear, 109, 189, 222
 polynomial, 108
 simple linear, 105–106, 108, 206
 reliability, 8
 replacement
 drawing with or without —, 19
 resampling, 203, 205
 residuals, 127, 147, 288
 response surface model, 109
 response variable, 105
 reversibility, 212
 R^2 , *see* coefficient of determination

S
 sample
 correlation coefficient, 124, 125, 156
 mean, 122, 123, 124
 median, 204
 standard deviation, 124
 variance, 123, 124, 206
 pooled, 133
 sample space, 5
 continuous, 11
 discrete, 10
 Savage–Dickey density ratio, 254
 score
 efficient, 167
 function, 165, 167
 interval, 174, 175
 seed, 52
 simple linear regression, 105–106, 115, 127, 139
 sort.m, 21
 sparse matrix, 295, 299, 307, 329, 362
 spreadsheet, 115
 standard deviation, 32
 sample, 124
 standard normal distribution, 46
 state space model, 323
 initial condition, 326
 stationarity, 289, 291
 statistic, 122, 140
 sufficient, *see* sufficient statistic
 statistical model, 102
 statistical test
 goodness of fit, 220
 steps for, 129, 141
 statistics, 3, 5
 Bayesian, 121
 classical, 121

stochastic volatility model, 339–345
 Student's t distribution, 50, 89, 131, 133, 266
 multivariate, 271, 285
 sufficient statistic, 150, 151, 153, 188
 sum rule, 9, 10, 16, 26, 27, 64, 65

T

target distribution, 209
 Taylor's theorem, 91
 multidimensional, 92, 108, 176, 177, 179, 180, 369
 test statistic, 140
 time series, 287–305, 323–345
 time-varying parameter autoregressive model, 333–339
 tower property, 78
 transformation
 of data, 110
 transformation rule, 79, 81, 242
 transition
 density, 210
 equation, 323
 graph, 210
 trimmed mean, 221
 truncated normal distribution, 279, 286
 two-sample
 binomial model, 103, 136
 normal model, 103, 111, 133

U

unbiased estimator, 122
 uniform distribution, 42, 188
 discrete, 59
 unobserved components model, 325–333

V

variance, 32
 properties, 33, 35, 36, 76, 77, 94
 sample, 123, 124, 206

W

Weibull
 distribution, 61, 190, 200
 model, 109

Z

zero inflated Poisson, 258