

PREFACE

We believe that many foundational ideas of Probability and Statistics are best understood when their natural connection is emphasised. We feel that the interested student should learn the mathematical rigour of Probability, the motivating examples and techniques from Statistics, and an instructive technology to perform computations relating to both in an inclusive manner. These formed our main motivations for writing this book. We have chosen to use the R software environment to demonstrate an available computational tool.

The book is intended to be an undergraduate text for a course on Probability Theory. We had in mind courses such as the one year (two semester) Probability course at many universities in India such as the Indian Statistical Institute or Chennai Mathematical Institute, or a one semester (or two quarter) Probability course as is commonly offered as an upper division, post-calculus elective at many North American universities. The Statistics material and the package R are introduced so as to emphasise motivations and applications of the probabilistic material. We assume that our readers are well-versed in calculus, have a basic understanding of the theory of sets and functions, combinatorics, and proof techniques, and have at least a passing awareness of the distinction between countable and uncountable infinities. We do not assume any particular experience of Linear Algebra or Real Analysis.

In Chapter 1 of this book we begin with an introduction to Outcomes, Sample Space, Events and the axiomatic definition of Probability. Then we discuss the concepts of conditional probability, independence and Bayes' Theorem. We conclude this chapter with a basic introduction to R. R is a Free Open Source software environment that runs on all major software platforms, and instructions to download and install it are available at <https://www.r-project.org/>.

We begin Chapter 2 by applying the notion of independence to repeated trials (Bernoulli Trials) and discuss the Binomial and Geometric distributions. We introduce the Poisson distribution as a limiting approximation of the Binomial. We conclude this section with a discussion on Sampling with and without replacement. The Hypergeometric Distribution is thus introduced here and we prove its approximation to the Binomial. Throughout this chapter and later in the book we provide the R code for calculating the probabilities associated with common distributions.

In Chapters 3 and 4 we introduce discrete random variables (functions on a sample space whose range is countable) and related concepts. In Chapter 3, we define the probability mass function, distribution function, and independence for random variables. We introduce the Multinomial distribution and show the memoryless property of the Geometric random variable. The chapter concludes by providing a method to compute the distributions of functions of one and several random variables, defining the concept of joint distribution along the way. In Chapter 4, we define Expectation, Variance, Covariance, Conditional Expectation and Conditional Variance for discrete random variables. Results involving these quantities for standard distributions are presented (and proved) as well. We also state and prove the Markov and Chebyshev inequalities along with the notion of standardising random variables to mean zero and variance one.

Working with uncountable spaces and understanding the probability density function of an absolutely continuous random variable are challenging without assuming a background in Real Analysis but we make a modest attempt towards this in Chapter 5. We begin with a description of uncountable sample spaces. After having described events in a temporary manner in Chapter 1 we provide a precise definition here but comfort the reader that we shall avoid the most general events and at most consider countable union/intersection of intervals. This allows us to be fairly rigorous

with random variables having piecewise continuous probability density functions using results from basic calculus. After this we imitate the program conducted in Chapter 3. Standard distributions such as Uniform, Exponential, and Normal are discussed. While computing densities of sums and ratios of independent random variables we introduce the Gamma distribution and use it to derive the Beta distribution as an example of ratio of dependent Gamma random variables.

In Chapter 6 we define Variance, Covariance, Conditional Expectation and Conditional Variance for continuous random variables and summarise their properties. Moment generating functions for random variables are defined. At this point, to respect the minimal background assumption on our reader we state a few important results without proof such as the fact that the moment generating functions characterise distribution of a random variable. The chapter ends with a section on Bivariate Normal random variables. Here we have done all computations in this section without using Linear Algebra but the notational efficiency of using Linear Algebra is explained via exercises for the interested reader.

With the foundational ideas of Probability laid out we proceed in Chapter 7 with Sampling and Descriptive Statistics. The empirical distribution, the sample mean, variance and proportion are defined along with their properties. Simulation is used to develop intuition regarding sampling variability, and plots such as Histograms, Hanging Rootograms, and Q-Q Plots are introduced and illustrated using R.

Limit Theorems for Sampling Distributions discussed in Chapter 7 are the objective of Chapter 8. We begin with a brief description of multivariate joint densities and Order statistics. The t -distribution and χ^2 (chi-square) distributions are introduced in this chapter. The sample mean and variance from a normal population are discussed in relation to t and χ^2 . We prove the Weak Law of Large numbers and the Central Limit Theorem for random variables possessing a moment generating function. We do state a more general version of the Central Limit Theorem and also the Strong Law of Large numbers, providing a proof of the latter in the Appendix. Along with R code we discuss the continuity correction and applications of the Central Limit Theorem via examples.

We end the book with two chapters focused solely on results and techniques from statistics. In Chapter 9 we discuss Estimation, Confidence Intervals and Hypothesis Testing. We briefly describe Method of Moments estimators followed by Maximum Likelihood Estimation. We then introduce Confidence Intervals as a methodology to work with when a single estimate might not suffice for the mean. In Hypothesis Testing we focus on tests involving Normal, t , χ^2 , and F distributions. We describe the use of traditional lookup tables to compute probabilities as well as the use of R as a more flexible alternative. We also discuss critical values and the concept of rejection region. In Chapter 10 we discuss elementary results from simple linear regression involving one independent variable and one dependent variable. We discuss the Least Squares method and also prove results that will be useful for predicting new data from the given model. In the final section we discuss applications of Hypothesis testing in Regression. As in earlier chapters, computations are illustrated using R.

R code for most of the computations done are given in the book itself, and the reader should be able to reproduce and extend them easily. Code for figures are not given in the book, but are available at a website accompanying the book, which also contains additional material for readers who are interested in learning more about R.

The Appendix includes proofs of results such as Strong Law of Large Numbers and other proofs that are of interest but beyond the scope of the text.

Siva Athreya, Deepayan Sarkar, and Steve Tanner
April 25, 2016

Probability and Statistics with Examples using R

Siva Athreya, Deepayan Sarkar, and Steve Tanner

April 25, 2016

CONTENTS

Preface	v
1 BASIC CONCEPTS	1
1.1 Definitions and Properties	1
1.1.1 Definitions	1
1.1.2 Basic Properties	2
1.2 Equally Likely Outcomes	8
1.3 Conditional Probability and Bayes' Theorem	12
1.3.1 Bayes' Theorem	17
1.4 Independence	20
1.5 Using R for computation	24
2 SAMPLING AND REPEATED TRIALS	29
2.1 Bernoulli Trials	29
2.1.1 Using R to compute probabilities	34
2.2 Poisson Approximation	37
2.3 Sampling With and Without Replacement	42
2.3.1 The Hypergeometric Distribution	43
2.3.2 Hypergeometric Distributions as a Series of Dependent Trials	43
2.3.3 Binomial Approximation to the Hypergeometric Distribution	45
3 DISCRETE RANDOM VARIABLES	49
3.1 Random Variables as Functions	49
3.1.1 Common Distributions	52
3.2 Independent and Dependent Variables	54
3.2.1 Independent Variables	55
3.2.2 Conditional, Joint, and Marginal Distributions	56
3.2.3 Memoryless Property of the Geometric Random Variable	59
3.2.4 Multinomial Distributions	60
3.3 Functions of Random Variables	63
3.3.1 Distribution of $f(X)$ and $f(X_1, X_2, \dots, X_n)$	63
3.3.2 Functions and Independence	67
4 SUMMARIZING DISCRETE RANDOM VARIABLES	71
4.1 Expected Value	71
4.1.1 Properties of the Expected Value	73
4.1.2 Expected Value of a Product	75
4.1.3 Expected Values of Common Distributions	76
4.1.4 Expected Value of $f(X_1, X_2, \dots, X_n)$	80
4.2 Variance and Standard Deviation	84
4.2.1 Properties of Variance and Standard Deviation	85
4.2.2 Variances of Common Distributions	87
4.2.3 Standardized Variables	89
4.3 Standard Units	93

4.3.1	Markov and Chebyshev Inequalities	94
4.4	Conditional Expectation and Conditional Variance	97
4.5	Covariance and Correlation	102
4.5.1	Covariance	103
4.5.2	Correlation	105
4.6	Exchangeable Random Variables	107
5	CONTINUOUS PROBABILITIES AND RANDOM VARIABLES	111
5.1	Uncountable Sample Spaces and Densities	111
5.1.1	Probability Densities on \mathbb{R}	113
5.2	Continuous Random Variables	117
5.2.1	Common Distributions	119
5.2.2	A word about individual outcomes	125
5.3	Transformation of Continuous Random Variables	130
5.4	Multiple Continuous Random Variables	137
5.4.1	Marginal Distributions	140
5.4.2	Independence	141
5.4.3	Conditional Density	143
5.5	Functions of Independent Random variables	148
5.5.1	Distributions of Sums of Independent Random variables	149
5.5.2	Distributions of Quotients of Independent Random variables	153
6	SUMMARISING CONTINUOUS RANDOM VARIABLES	161
6.1	Expectation, and Variance	161
6.2	Covariance, Correlation, Conditional Expectation and Conditional Variance	168
6.3	Moment Generating Functions	176
6.4	Bivariate Normals	181
7	SAMPLING AND DESCRIPTIVE STATISTICS	187
7.1	The empirical distribution	187
7.2	Descriptive Statistics	188
7.2.1	Sample Mean	188
7.2.2	Sample Variance	189
7.2.3	Sample proportion	189
7.3	Simulation	191
7.4	Plots	194
7.4.1	Empirical Distribution Plot for Discrete Distributions	195
7.4.2	Histograms for Continuous Distributions	197
7.4.3	Hanging Rootograms for Comparing with Theoretical Distributions	197
7.4.4	Q-Q Plots for Continuous Distributions	200
8	SAMPLING DISTRIBUTIONS AND LIMIT THEOREMS	203
8.1	Multi-dimensional continuous random variables	203
8.1.1	Order Statistics and their Distributions	205
8.1.2	χ^2 , F and t	208
8.1.3	Distribution of Sampling Statistics from a Normal population	210
8.2	Weak Law of Large Numbers	214
8.3	Convergence in Distribution	216
8.4	Central Limit Theorem	218

8.4.1	Normal Approximation and Continuity Correction	220
9	ESTIMATION AND HYPOTHESIS TESTING	225
9.1	Notations and Terminology for Estimators	225
9.2	Method of Moments	226
9.3	Maximum Likelihood Estimate	227
9.4	Confidence Intervals	229
9.4.1	Confidence Intervals when the standard deviation σ is known	229
9.4.2	Confidence Intervals when the standard deviation σ is unknown	230
9.5	Hypothesis Testing	231
9.5.1	The z-test: Test for sample mean when σ is known	231
9.5.2	The t-test: Test for sample mean when σ is unknown	233
9.5.3	A critical value approach	234
9.5.4	The χ^2 -test : Test for sample variance	234
9.5.5	The two-sample z-test: Test to compare sample means	236
9.5.6	The F-test: Test to compare sample variances.	237
9.5.7	A χ^2 -test for “goodness of fit”	237
10	LINEAR REGRESSION	241
10.1	Sample Covariance and Correlation	241
10.2	Simple Linear Model	241
10.3	The Least Squares Line	242
10.4	a and b as Random Variables	244
10.5	Predicting New Data When σ^2 is Known	247
10.6	Hypothesis Testing and Regression	249
10.7	Estimaing an Unknown σ^2	249
A	WORKING WITH DATA IN R	253
A.1	Datasets in R	253
A.2	Plotting data	254
B	SOME MATHEMATICAL DETAILS	255
B.1	Linear Algebra	255
B.2	Jacobian Method	255
B.3	χ^2 -goodness of fit test	255
C	STRONG LAW OF LARGE NUMBERS	257
D	TABLES	261

I

BASIC CONCEPTS

1.1 DEFINITIONS AND PROPERTIES

Most of the problems in probability and statistics involve determining how likely it is that certain things will occur. Before we can talk about what is likely or unlikely, we need to know what is possible. In other words, we need some framework in which to discuss what sorts of things have the potential to occur. To that end, we begin by introducing the basic concepts of “sample space”, “experiment”, “outcome”, and “event”. We also define what we mean by a “probability” and provide some examples to demonstrate the consequences of the definition.

1.1.1 Definitions

DEFINITION 1.1.1. (Sample Space) *A sample space S is a set. The elements of the set S will be called “outcomes” and should be viewed as a listing of all possibilities that might occur. We will call the process of actually selecting one of these outcomes an “experiment”.*

For its familiarity and simplicity we will frequently use the example of rolling a die. In that case our sample space would be $S = \{1, 2, 3, 4, 5, 6\}$, a complete listing of all of the outcomes on the die. Performing an experiment in this case would mean rolling the die and recording the number that it shows. However, sample space outcomes need not be numeric. If we are flipping a coin (another simple and familiar example) experiments would result in one of two outcomes and the appropriate sample space would be $S = \{\text{Heads}, \text{Tails}\}$.

For a more interesting example, if we are discussing which country will win the next World Cup, outcomes might include Brazil, Spain, Canada, and Thailand. Here the set S might be all the world’s countries. An experiment in this case requires waiting for the next World Cup and identifying the country which wins the championship game. Though we have not yet explained how probability relates to a sample space, soccer fans amongst our readers may regard this example as a foreshadowing that not all outcomes of a sample space will necessarily have the same associated probabilities.

DEFINITION 1.1.2. (Temporary Definition of Event) *Given a sample space S , an “event” is any subset $E \subset S$.*

This definition will allow us to talk about how likely it is that a range of possible outcomes might occur. Continuing our examples above we might want to talk about the probability that a die rolls a number larger than two. This would involve the event $\{3, 4, 5, 6\}$ as a subset of $\{1, 2, 3, 4, 5, 6\}$. In the soccer example we might ask whether the World Cup will be won by a South American

country. This subset of our list of all the world's nations would contain Brazil as an element, but not Spain.

It is worth noting that the definition of “event” includes both S , the sample space itself, and \emptyset , the empty set, as legitimate examples. As we introduce more complicated examples we will see that it is not always necessary (or even possible) to regard every single subset of a sample space as a legitimate event, but since the reasons for that may be distracting at this point we will use the above as a temporary definition of “event” and refine the definition when it becomes necessary.

To each event, we want to assign a chance (or “probability”) which will be a number between 0 and 1. So if the probability of an event E is 0.72, we interpret that as saying, “When an experiment is performed, it has a 72% chance of resulting in an outcome contained in the event E ”. Probabilities will satisfy two axioms stated and explained below. This formal definition is due to Andrey Kolmogorov (1903-1987).

DEFINITION 1.1.3. (Probability Space Axioms) *Let S be a sample space and let \mathcal{F} be the collection of all events.*

A “probability” is a function $P : \mathcal{F} \rightarrow [0, 1]$ such that

- (1) $P(S) = 1$; and
- (2) If E_1, E_2, \dots are a countable collection of disjoint events
(that is, $E_i \cap E_j = \emptyset$ if $i \neq j$), then

$$P\left(\bigcup_{j=1}^{\infty} E_j\right) = \sum_{j=1}^{\infty} P(E_j). \quad (1.1.1)$$

The first axiom is relatively straight forward. It simply reiterates that S did, indeed, include all possibilities, and therefore there is a 100% chance that an experiment will result in some outcome included in S . The second axiom is not as complicated as it looks. It simply says that probabilities add when combining a countable number of disjoint events. It is implicit that the series on right hand side of the equation (1.1.1) converges. Further (1.1.1) also holds when combining finite number of disjoint events (see Theorem 1.1.4 below).

Returning to our World Cup example, suppose A is a list of all North American countries and E is a list of all European countries. If it happens that $P(A) = 0.05$ and $P(E) = 0.57$ then $P(A \cup E) = 0.62$. In other words, if there is a 5% chance the next World Cup will be won by a North American nation and a 57% chance that it will be won by a European nation, then there is a 62% chance that it will be won by a nation from either Europe or North America. The disjointness of these events is obvious since (if we discount island territories) there isn't any country that is in both North America and Europe.

The requirement of axiom two that the collection of events be countable is important. We shall see shortly that, as a consequence of axiom two, disjoint additivity also applies to any finite collection of events. It does not apply to uncountably infinite collections of events, though that fact will not be relevant until later in the text when we discuss continuous probability spaces.

1.1.2 Basic Properties

There are some immediate consequences of these probability axioms which we will state and prove before returning to some simple examples.

THEOREM 1.1.4. *Let P be a probability on a sample space S . Then,*

- (1) $P(\emptyset) = 0$;
- (2) *If E_1, E_2, \dots, E_n are a finite collection of disjoint events, then*

$$P\left(\bigcup_{j=1}^n E_j\right) = \sum_{j=1}^n P(E_j);$$

- (3) *If E and F are events with $E \subset F$, then $P(E) \leq P(F)$;*
- (4) *If E and F are events with $E \subset F$, then $P(F \setminus E) = P(F) - P(E)$;*
- (5) *Let E^c be the complement of event E . Then $P(E^c) = 1 - P(E)$; and*
- (6) *If E and F are events then $P(E \cup F) = P(E) + P(F) - P(E \cap F)$.*

Proof of (1) - The empty set is disjoint from itself, so $\emptyset, \emptyset, \dots$ is a countable disjoint collection of events. From the second axiom, $P\left(\bigcup_{j=1}^{\infty} E_j\right) = \sum_{j=1}^{\infty} P(E_j)$. When this is applied to the collection of empty sets we have $P(\emptyset) = \sum_{j=1}^{\infty} P(\emptyset)$. If $P(\emptyset)$ had any non-zero value, the right hand side of this equation would be a divergent series while the left hand side would be a number. Therefore, $P(\emptyset) = 0$.

Proof of (2) - To use axiom two we need to make this a countable collection of events. We may do so while preserving disjointness by including copies of the empty set. Define $E_j = \emptyset$ for $j > n$. Then $E_1, E_2, \dots, E_n, \emptyset, \emptyset, \dots$ is a countable collection of disjoint sets and therefore $P\left(\bigcup_{j=1}^{\infty} E_j\right) = \sum_{j=1}^{\infty} P(E_j)$. However, the empty sets add nothing to the union and so $\bigcup_{j=1}^{\infty} E_j = \bigcup_{j=1}^n E_j$. Likewise since we have shown $P(\emptyset) = 0$ these sets also add nothing to the sum, so $\sum_{j=1}^{\infty} P(E_j) = \sum_{j=1}^n P(E_j)$.

Combining these gives the result:

$$P\left(\bigcup_{j=1}^n E_j\right) = P\left(\bigcup_{j=1}^{\infty} E_j\right) = \sum_{j=1}^{\infty} P(E_j) = \sum_{j=1}^n P(E_j).$$

Proof of (3) - If $E \subset F$, then E and $F \setminus E$ are disjoint events with a union equal to F . Using (2) above gives $P(F) = P(E \cup (F \setminus E)) = P(E) + P(F \setminus E)$.

Since probabilities are assumed to be positive, it follows that $P(F) \geq P(E)$.

Proof of (4) - As with the proof of (3) above, E and $F \setminus E$ are disjoint events with $E \cup (F \setminus E) = F$. Therefore $P(F) = P(E) + P(F \setminus E)$ from which we get the result.

Proof of (5) - This is simple a special case of (4) where $F = S$.

Proof of (6) - We may disassemble $E \cup F$ disjointly as $E \cup F = E \cup (F \setminus E)$. Then from (2) we have $P(E \cup F) = P(E) + P(F \setminus E)$.

Next, since $F \setminus E \subset F$ and since $F \setminus (F \setminus E) = E \cap F$ we can use (4) to write $P(E \cup F) = P(E) + P(F) - P(E \cap F)$. ■

EXAMPLE 1.1.5. A coin flip can come up either “heads” or “tails”, so $S = \{\text{heads}, \text{tails}\}$. A coin is considered “fair” if each of these outcomes is equally likely. Which axioms or properties above can be used to reach the (obvious) conclusion that both outcomes have a 50% chance of occurring?

Each outcome can also be regarded as an event. So $E = \{\text{heads}\}$ and $F = \{\text{tails}\}$ are two disjoint events. If the coin is fair, each of these events is equally likely, so $P(E) = P(F) = p$ for some value of p . However, using the second axiom, $1 = P(S) = P(E \cup F) = P(E) + P(F) = 2p$. Therefore, $p = 0.5$, or in other words each of the two possibilities has a 50% chance of occurring on any flip of a fair coin. ■

In the examples above we have explicitly described the sample space S , but in many cases this is neither necessary nor desirable. We may still use the probability space axioms and their consequences when we know the probabilities of certain events even if the sample space is not explicitly described.

EXAMPLE 1.1.6. A certain sea-side town has a small fishing industry. The quantity of fish caught by the town in a given year is variable, but we know there is a 35% chance that the town’s fleet will catch over 400 tons of fish, but only a 10% chance that they will catch over 500 tons of fish. How likely is it they will catch between 400 and 500 tons of fish?

The answer to this may be obvious without resorting to sets, but we use it as a first example to illustrate the proper use of events. Note, though, that we will not explicitly describe the sample space S .

There are two relevant events described in the problem above. We have F representing “the town’s fleet will catch over 400 tons of fish” and E representing “the town’s fleet will catch over 500 tons of fish”. We are given that $P(E) = 0.1$ while $P(F) = 0.35$.

Of course $E \subset F$ since if over 500 tons of fish are caught, the actual tonnage will be over 400 as well. The event that the town’s fleet will catch between 400 and 500 tons of fish is $F \setminus E$ since E hasn’t occurred, but F has. So using property (4) from above we have $P(F \setminus E) = P(F) - P(E) = 0.35 - 0.1 = 0.25$. In other words there is a 25% chance that between 400 and 500 tons of fish will be caught. ■

EXAMPLE 1.1.7. Suppose we know there is a 60% chance that it will rain tomorrow and a 70% chance the high temperature will be above 30°C . Suppose we also know that there is a 40% chance that the high temperature will be above 30°C and it will rain. How likely is it tomorrow will be a dry day that does not go above 30°C ?

The answer to this question may not be so obvious, but our first step is still to view the pieces of information in terms of events and probabilities. We have one event E which represents “It will rain tomorrow” and another F which represents “The high will be above 30°C tomorrow”. Our given probabilities tell us $P(E) = 0.6$, $P(F) = 0.7$, and $P(E \cap F) = 0.4$. We are trying to determine $P(E^c \cap F^c)$. We can do so using properties (5) and (6) proven above, together with the set-algebraic fact that $E^c \cap F^c = (E \cup F)^c$.

From (5) we know $P(E \cup F) = P(E) + P(F) - P(E \cap F) = 0.7 + 0.6 - 0.4 = 0.9$. (This is the probability that it either will rain or be above 30°C).

Then from (6) and the set-algebraic fact, $P(E^c \cap F^c) = P((E \cup F)^c) = 1 - P(E \cup F) = 1 - 0.9 = 0.1$.

So there is a 10% chance tomorrow will be a dry day that doesn’t reach 30 degrees. ■

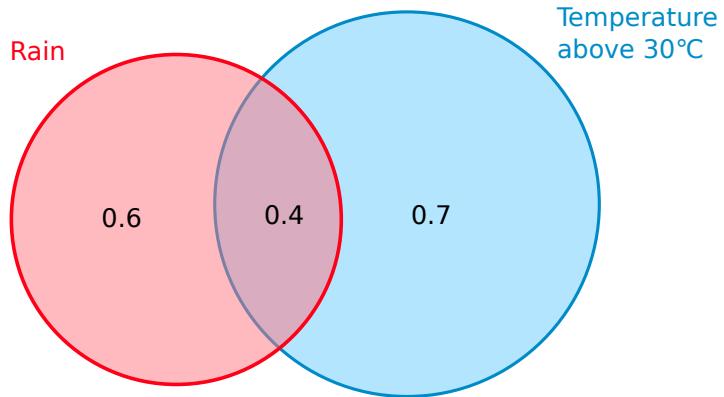


Figure 1.1: A Venn diagram that describes the probabilities from Example 1.1.7.

EXERCISES

Ex. 1.1.1. Consider the sample space $\Omega = \{a, b, c, d, e\}$. Given that $\{a, b, e\}$, and $\{b, c\}$ are both **events**, what other subsets of Ω must be events due to the requirement that the collection of events is closed under taking unions, intersections, and compliments?

Ex. 1.1.2. There are two positions - Cashier and Waiter - open at the local restaurant. There are two male applicants (David and Rajesh) two female applicants (Veronica and Megha). The Cashier position is chosen by selecting one of the four applicants at random. The Waiter position is then chosen by selecting at random one of the three remaining applicants.

- (a) List the elements of the sample space S .
- (b) List the elements of the event A that the position of Cashier is filled by a female applicant.
- (c) List the elements of the event B that exactly one of the two positions is filled by a female applicant.
- (d) List the elements of the event C that neither position was filled by a female applicant.
- (e) Sketch a Venn diagram to show the relationship among the events A , B , C and S .

Ex. 1.1.3. A jar contains a large collection of red, green, and white marbles. Marbles are drawn from the jar one at a time. The color of the marble is recorded and it is put back in the jar before the next draw. Let R_n denote the event that the n^{th} draw is a red marble and let G_n denote the event that the n^{th} draw is a green marble. For example, $R_1 \cap G_2$ would denote the event that the first marble was red and the second was green. In terms of these events (and appropriate set-theoretic symbols – union, intersection, and complement) find expressions for the events in parts (a), (b), and (c) below.

- (a) The first marble drawn is white. (We might call this W_1 , but show that it can be written in terms of the R_n and G_n sets described above).

- (b) The first marble drawn is green and the second marble drawn is not white.
- (c) The first and second draws are different colors.
- (d) Let $E = R_1 \cup G_2$ and let $F = R_1^c \cap R_2$. Are E and F disjoint? Why or why not?

Ex. 1.1.4. Suppose there are only thirteen teams with a non-zero chance of winning the next World Cup. Suppose those teams are Spain (with a 14% chance), the Netherlands (with a 11% chance), Germany (with a 11% chance), Italy (with a 10% chance), Brazil (with a 10% chance), England (with a 9% chance), Argentina (with a 9% chance), Russia (with a 7% chance), France (with an 6% chance), Turkey (with a 4% chance), Paraguay (with a 4% chance), Croatia (with a 4% chance) and Portugal (with a 1% chance).

- (a) What is the probability that the next World Cup will be won by a South American country?
- (b) What is the probability that the next World Cup will be won by a country that is not from South America? (Think of two ways to do this problem - one directly and one using part (5) of Theorem 1.1.4. Which do you prefer and why?)

Ex. 1.1.5. If A and B are disjoint events and $P(A) = 0.3$ and $P(B) = 0.6$, find $P(A \cup B)$, $P(A^c)$ and $P(A^c \cap B)$.

Ex. 1.1.6. Suppose E and F are events in a sample space S . Suppose that $P(E) = 0.7$ and $P(F) = 0.5$.

- (a) What is the largest possible value of $P(E \cap F)$? Explain.
- (b) What is the smallest possible value of $P(E \cap F)$? Explain.

Ex. 1.1.7. A biologist is modeling the size of a frog population in a series of ponds. She is concerned with both the number of egg masses laid by the frogs during breeding season and the annual precipitation into the ponds. She knows that in a given year there is an 86% chance that there will be over 150 egg masses deposited by the frogs (event E) and that there is a 64% chance that the annual precipitation will be over 17 inches (event F).

- (a) In terms of E and F , what is the event “there will be over 150 egg masses and an annual precipitation of over 17 inches”?
- (b) In terms of E and F , what is the event “there will be 150 or fewer egg masses and the annual precipitation will be over 17 inches”?
- (c) Suppose the probability of the event from (a) is 59%. What is the probability of the event from (b)?

Ex. 1.1.8. In part (6) of Theorem 1.1.4 we showed that

$$P(E \cup F) = P(E) + P(F) - P(E \cap F).$$

Versions of this rule for three or more sets are explored below.

- (a) Prove that $P(A \cup B \cup C)$ is equal to

$$P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$$

for any events A , B , and C .

- (b) Use part (a) to answer the following question. Suppose that in a certain United States city 49.3% of the population is male, 11.6% of the population is sixty-five years of age or older, and 13.8% of the population is Hispanic. Further, suppose 5.1% is both male and at least sixty-five, 1.8% is both male and Hispanic, and 5.9% is Hispanic and at least sixty-five. Finally, suppose that 0.7% of the population consists of Hispanic men that are at least sixty-five years old. What percentage of people in this city consists of non-Hispanic women younger than sixty-five years old?
- (c) Find a four-set version of the equation. That is, write $P(A \cup B \cup C \cup D)$ in terms of probabilities of intersections of A, B, C , and D .
- (d) Find an n-set version of the equation.

Ex. 1.1.9. A and B are two events. $P(A)=0.4$, $P(B)=0.3$, $P(A \cup B)=0.6$. Find the following probabilities:

- (a) $P(A \cap B)$;
- (b) $P(\text{Only } A \text{ happens})$; and
- (c) $P(\text{Exactly one of } A \text{ or } B \text{ happens})$.

Ex. 1.1.10. In the next subsection we begin to look at probability spaces where each of the outcomes are equally likely. This problem will help develop some early intuition for such problems.

- (a) Suppose we roll a die and so $S = \{1, 2, 3, 4, 5, 6\}$. Each outcome separately $\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}$ is an event. Suppose each of these events is equally likely. What must the probability of each event be? What axioms or properties are you using to come to your conclusion?
- (b) With the same assumptions as in part (a), how would you determine the probability of an event like $E = \{1, 3, 4, 6\}$? What axioms or properties are you using to come to your conclusion?
- (c) If $S = \{1, 2, 3, \dots, n\}$ and each single-outcome event is equally likely, what would be the probability of each of these events?
- (d) Suppose $E \subset S$ is an event in the sample space from part (c). Explain how you could determine $P(E)$.

Ex. 1.1.11. Suppose A and B are subsets of a sample space Ω .

- (a) Show that $(A - B) \cup B = A$ when $B \subset A$.
- (b) Show by example that the equality doesn't always hold if B is not a subset of A .

Ex. 1.1.12. Let A and B be events.

- (a) Suppose $P(A) = P(B) = 0$. Prove that $P(A \cup B) = 0$.
- (b) Suppose $P(A) = P(B) = 1$. Prove that $P(A \cap B) = 1$.

Ex. 1.1.13. Let A_n be a sequence of events.

- (a) Suppose $A_n \subseteq A_{n+1}$ for all $n \geq 1$. Show that

$$P(\bigcup_{n=1}^{\infty} A_n) = \lim_{n \rightarrow \infty} P(A_n)$$

- (b) Suppose $A_n \supseteq A_{n+1}$ for all $n \geq 1$. Show that

$$P(\bigcap_{n=1}^{\infty} A_n) = \lim_{n \rightarrow \infty} P(A_n)$$

1.2 EQUALLY LIKELY OUTCOMES

When a sample space S consists of only a countable collection of outcomes, describing the probability of each individual outcome is sufficient to describe the probability of all events. This is because if $A \subset S$ we may simply compute

$$P(A) = P\left(\bigcup_{\omega \in A} \{\omega\}\right) = \sum_{\omega \in A} P(\{\omega\}).$$

This assignment of probabilities to each outcome is called a “distribution” since it describes how probability is distributed amongst the possibilities. Perhaps the simplest example arises when there are a finite collection of equally likely outcomes. Think of examples such as flipping a fair coin (“heads” and “tails” are equally likely to occur), rolling a fair die (1, 2, 3, 4, 5, and 6 are equally likely), or drawing a set of numbers for a lottery (many possibilities, but in a typical lottery, each outcome is as likely as any other). Such distributions are common enough that it is useful to have shorthand notations for them. In the case of a sample space $S = \{\omega_1, \dots, \omega_n\}$ where each outcome is equally likely, the probability is referred to as a “uniform distribution” and is denoted by $\text{Uniform}(\{\omega_1, \dots, \omega_n\})$. In such situations, computing probabilities simply reduces to computing the number of outcomes in a given event and consequently becomes a combinatorial problem.

THEOREM 1.2.1. Uniform($\{\omega_1, \dots, \omega_n\}$) : Let $S = \{\omega_1, \omega_2, \dots, \omega_n\}$ be a non-empty, finite set. If $E \subset S$ is any subset of S , let $P(E) = \frac{|E|}{|S|}$ (where $|E|$ represents the number of elements of E). Then P defines a probability on S and P assigns equal probability to each individual outcome in S .

Proof - Since $E \subset S$ we know $|E| \leq |S|$ and so $0 \leq P(E) \leq 1$, so we must prove that P satisfies the two probability axioms.

Since $P(S) = \frac{|S|}{|S|} = 1$ the first axiom is satisfied.

To verify the second axiom, suppose E_1, E_2, \dots is a countable collection of disjoint events. Since S is finite, only finitely many of these E_j can be non-empty, so we may list the non-empty events as E_1, E_2, \dots, E_n . For $j > n$ we know $E_j = \emptyset$ and so $P(E_j) = 0$ by the definition. Since the events are disjoint, to find the number of elements in their union we simply add the elements of each event separately. That is, $|E_1 \cup E_2 \cup \dots \cup E_n| = |E_1| + |E_2| + \dots + |E_n|$ and therefore

$$P\left(\bigcup_{j=1}^{\infty} E_j\right) = \frac{|\bigcup_{j=1}^{\infty} E_j|}{|S|} = \frac{\sum_{j=1}^n |E_j|}{|S|} = \sum_{j=1}^n \frac{|E_j|}{|S|} = \sum_{j=1}^n P(E_j) = \sum_{j=1}^{\infty} P(E_j).$$

Finally, let $\omega \in S$ be any single outcome and let $E = \{\omega\}$. Then $P(E) = \frac{1}{|S|}$, so every outcome in S is equally likely. ■

EXAMPLE 1.2.2. A deck of twenty cards labeled 1, 2, 3, ..., 20 is shuffled and a card selected at random. What is the probability that the number on the card is a multiple of six?

The description of the scenario suggests that each of the twenty cards is as likely to be chosen as any other. In this case $S = \{1, 2, 3, \dots, 20\}$ while $E = \{6, 12, 18\}$. Therefore, $P(E) = \frac{|E|}{|S|} = \frac{3}{20} = 0.15$. There is a 15% chance that the card will be a multiple of six. ■

EXAMPLE 1.2.3. Two dice are rolled. How likely is it that their sum will equal eight?

Since we are looking at a sum of dice, it might be tempting to regard the sample space as $S = \{2, 3, 4, \dots, 11, 12\}$, the collection of possible sums. While this is a possible approach (and one

we will return to later), it is not the case that all of these outcomes are equally likely. Instead we can view an experiment as tossing a first die and a second die and recording the pair of numbers that occur on each of the dice. Each of these pairs is as likely as any other to occur. So

$$S = \left\{ (1,1), (1,2), (1,3), (1,4), (1,5), (1,6), (2,1), (2,2), (2,3), (2,4), (2,5), (2,6), (3,1), (3,2), (3,3), (3,4), (3,5), (3,6), (4,1), (4,2), (4,3), (4,4), (4,5), (4,6), (5,1), (5,2), (5,3), (5,4), (5,5), (5,6), (6,1), (6,2), (6,3), (6,4), (6,5), (6,6) \right\}$$

and $|S| = 6 \times 6 = 36$. The event that the sum of the dice is an eight is $E = \{(2,6), (3,5), (4,4), (5,3), (6,2)\}$. Therefore $P(E) = \frac{|E|}{|S|} = \frac{5}{36}$. ■

EXAMPLE 1.2.4. A seven letter code is selected at random with every code as likely to be selected as any other code (so AQRVTAS and CRXAOLZ would be two possibilities). How likely is it that such a code has at least one letter used more than once? (This would happen with the first code above with a repeated A - but not with the second).

As with the examples above, the solution amounts to counting numbers of outcomes. However, unlike the examples above the numbers involved here are quite large and we will need to use some combinatorics to find the solution. The sample space S consists of all seven-letter codes from AAAAAAA to ZZZZZZZ. Each of the seven spots in the code could be any of twenty-six letters, so $|S| = 26^7 = 8,031,810,176$. If E is the event for which there is at least one letter used more than once, it is easier to count E^c , the event where no letter is repeated. Since in this case each new letter rules out a possibility for the next letter there are $26 \times 25 \times 24 \times 23 \times 22 \times 21 \times 20 = 3,315,312,000$ such possibilities.

This lets us compute $P(E^c) = \frac{3,315,312,000}{8,031,810,176}$ from which we find $P(E) = 1 - P(E^c) = \frac{4,716,498,176}{8,031,810,176} \approx 0.587$. That is, there is about a 58.7% chance that such a code will have a repeated letter. ■

EXAMPLE 1.2.5. A group of twelve people includes Grant and Dilip. A group of three people is to be randomly selected from the twelve. How likely is it that this three-person group will include Grant, but not Dilip?

Here, S is the collection of all three-person groups, each of which is as likely to be selected as any other. The number of ways of selecting a three-person group from a pool of twelve is $|S| = \binom{12}{3} = 220$. The event E consists of those three-person groups that include Grant, but not Dilip. Such groups must include two people other than Grant and there are ten people remaining from which to select the two, so $|E| = \binom{10}{2} = 45$. Therefore, $P(E) = \frac{45}{220} = \frac{9}{44}$. ■

EXERCISES

Ex. 1.2.1. A day is selected at random from a given week with each day as likely to be selected as any other.

- (a) What is the sample space S ? What is the size of S ?
- (b) Let E be the event that the selected day is a Saturday or a Sunday. What is the probability of E .

Ex. 1.2.2. A box contains 500 envelopes, of which 50 contain Rs 100 in cash, 100 contain Rs 50 in cash and 350 contain Rs 10. An envelope can be purchased at Rs 25 from the owner, who will pick

an envelope at random and give it to you. Write down the sample space for the net money gained by you. If each envelope is as likely to be selected as any other envelope, what is the probability that the first envelope purchased contains less than Rs 100?

Ex. 1.2.3. Three dice are tossed.

- (a) Describe (in words) the sample space S and give an example of an object in S .
- (b) What is the size of S ?
- (c) Let E be the event that the first two dice both come up “1”. What is the size of E ? What is the probability of E ?
- (d) Let G be the event that the three dice show three different numbers. What is the size of G ? What is the probability of G ?
- (e) Let F be the event that the third die is larger than the sum of the first two. What is the size of F ? What is the probability of F ?

Ex. 1.2.4. Suppose that each of three women at a party throws her hat into the center of the room. The hats are first mixed up and then each one randomly selects a hat. Describe the probability space for the possible selection of hats. If all of these selections are equally likely, what is the probability that none of the three women selects her own hat?

Ex. 1.2.5. A group of ten people includes Sona and Adam. A group of five people is to be randomly selected from the ten. How likely is it that this group of five people will include neither Sona nor Adam?

Ex. 1.2.6. There are eight students with two females and six males. They are split into two groups A and B, of four each.

- (a) In how many different ways can this be done?
- (b) What is the probability that two females end up in group A?
- (c) What is the probability that there is one female in each group?

Ex. 1.2.7. Sheela has lost her key to her room. The security officer gives her 50 keys and tells her that one of them will open her room. She decides to try each key successively and notes down the number of the attempt at which the room opens. Describe the sample space for this experiment. Do you think it is realistic that each of these outcomes is equally likely? Why or why not?

Ex. 1.2.8. Suppose that n balls, of which k are red, are arranged at random in a line. What is the probability that all the red balls are next to each other?

Ex. 1.2.9. Consider a deck of 50 cards. Each card has one of 5 colors (black, blue, green, red, and yellow), and is printed with a number (1,2,3,4,5,6,7,8,9, or 10) so that each of the 50 color/number combinations is represented exactly once. A **hand** is produced by dealing out five different cards from the deck. The order in which the cards were dealt does not matter.

- (a) How many different hands are there?
- (b) How many hands consist of cards of identical color? What is the probability of being dealt such a hand?
- (c) What is the probability of being dealt a hand that contains exactly three cards with one number, and two cards with a different number?

- (d) What is the probability of being dealt a hand that contains two cards with one number, two cards with a different number, and one card of a third number?

Ex. 1.2.10. Suppose you are in charge of quality control for a light bulb manufacturing company. Suppose that in the process of producing 100 light bulbs, either all 100 bulbs will work properly, or through some manufacturing error twenty of the 100 will not work. Suppose your quality control procedure is to randomly select ten bulbs from a 100 bulb batch and test them to see if they work properly. How likely is this procedure to detect if a batch has bad bulbs in it?

Ex. 1.2.11. A fair die is rolled five times. What is the probability of getting at least two 5's and at least two 6's among the five rolls.

Ex. 1.2.12. (**The “Birthday Problem”**) For a group of N people, if their birthdays were listed one-by-one, there are 365^N different ways that such a list might read (if we ignore February 29 as a possibility). Suppose each of those possible lists is as likely as any other.

- (a) For a group of two people, let E be the event that they have the same birthday. What is the size of E ? What is the probability of E ?
- (b) For a group of three people, let F be the event that at least two of the three have the same birthday. What is the size of F ? What is the probability of F ? (Hint: It is easier to find the size of F^c than it is to find the size of F).
- (c) For a group of four people, how likely is it that at least two of the four have the same birthday?
- (d) How large a group of people would you need to have before it becomes more likely than not that at least two of them share a birthday?

Ex. 1.2.13. A coin is tossed 100 times.

- (a) How likely is it that the 100 tosses will produce exactly fifty heads and fifty tails?
- (b) How likely is it that the number of heads will be between 50 and 55 (inclusive)?

Ex. 1.2.14. Suppose I have a coin that I claim is “fair” (equally likely to come up heads or tails) and that my friend claims is weighted towards heads. Suppose I flip the coin twenty times and find that it comes up heads on sixteen of those twenty flips. While this seems to favor my friend’s hypothesis, it is still possible that I am correct about the coin and that just by chance the coin happened to come up heads more often than tails on this series of flips. Let S be the sample space of all possible sequences of flips. The size of S is then 2^{20} , and if I am correct about the coin being “fair”, each of these outcomes are equally likely.

- (a) Let E be the event that exactly sixteen of the flips come up heads. What is the size of E ? What is the probability of E ?
- (b) Let F be the event that at least sixteen of the flips come up heads. What is the size of F ? What is the probability of F ?

Note that the probability of F is the chance of getting a result as extreme as the one I observed if I happen to be correct about the coin being fair. The larger $P(F)$ is, the more reasonable seems my assumption about the coin being fair. The smaller $P(F)$ is, the more that assumption looks doubtful. This is the basic idea behind the statistical concept of “hypothesis testing” which we will revisit in Chapter 9.

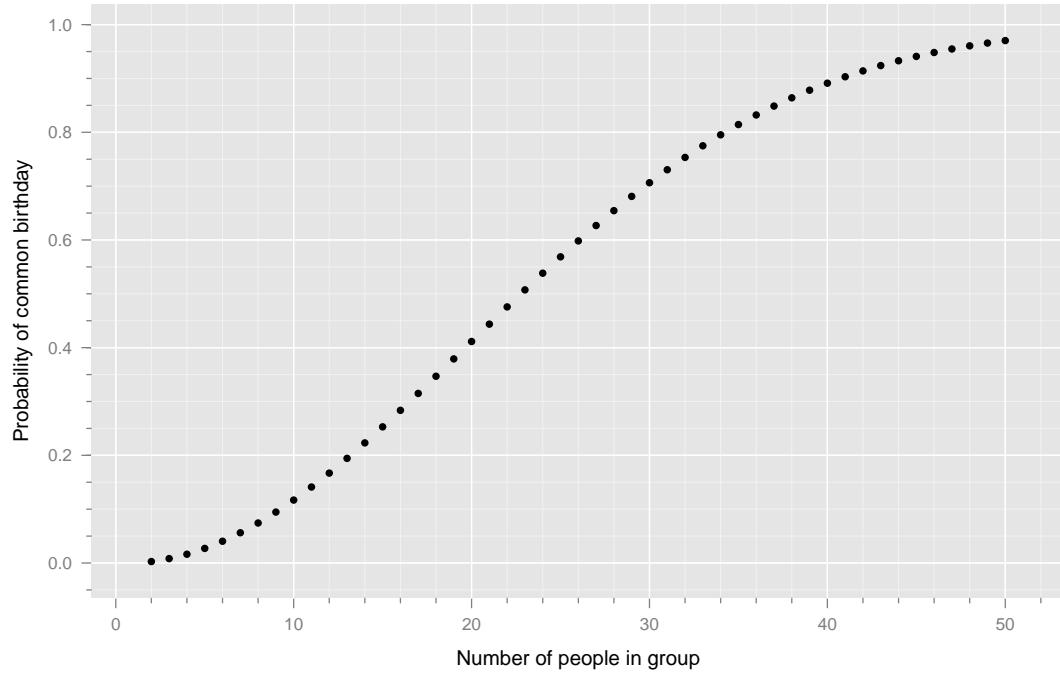


Figure 1.2: The birthday problem discussed in Exercise 1.2.12

Ex. 1.2.15. Suppose that r indistinguishable balls are placed in n distinguishable boxes so that each distinguishable arrangement is equally likely. Find the probability that no box will be empty.

Ex. 1.2.16. Suppose that 10 potato sticks are broken into two - one long and one short piece. The 20 pieces are now arranged into 10 random pairs chosen uniformly.

- (a) Find the probability that each of pairs consists of two pieces that were originally part of the same potato stick.
- (b) Find the probability that each pair consists of a long piece and a short piece.

Ex. 1.2.17. Let S be a non-empty, countable (finite or infinite) set such that for each $\omega \in S$, $0 \leq p_\omega \leq 1$. Let \mathcal{F} be the collection of all events. Suppose $P : \mathcal{F} \rightarrow [0, 1]$ is given by

$$P(E) = \sum_{\omega \in E} p_\omega,$$

for any event E .

- (a) Show that P satisfies Axiom 2 in Definition 1.1.3.
- (b) Further, conclude that if $P(S) = 1$ then P defines a probability on S .

1.3 CONDITIONAL PROBABILITY AND BAYES' THEOREM

In the previous section we introduced an axiomatic definition of “probability” and discussed the concept of an “event”. Now we look at ways in which the knowledge that one event has occurred may be used as information to inform and alter the probability of another event.

EXAMPLE 1.3.1. Consider the experiment of tossing a fair coin three times with sample space $S = \{hhh, hht, hth, htt, thh, tth, ttt\}$. Let A be the event that there are two or more heads. As all outcomes are equally likely,

$$\begin{aligned} P(A) &= \frac{|A|}{|S|} \\ &= \frac{|\{hhh, hht, hth, thh\}|}{8} \\ &= \frac{1}{2}. \end{aligned}$$

Let B be the event that there is a head in the first toss. As above,

$$\begin{aligned} P(B) &= \frac{|B|}{|S|} \\ &= \frac{|\{hhh, hht, hth, htt\}|}{8} \\ &= \frac{1}{2}. \end{aligned}$$

Now suppose we are asked to find the probability of at least two or more heads among the three tosses, but we are also given the additional information that the first toss was a head. In other words, we are asked to find the probability of A , given the information that event B has definitely occurred. Since the additional information guarantees B is now a list of all possible outcomes, it makes intuitive sense to view the event B as a new sample space and then identify the subset $A \cap B = \{hhh, hht, hth\}$ of B consisting of outcomes for which there are at least two heads. We could conclude that the probability of at least two or more heads in three tosses given that the first toss was a head is

$$\frac{|A \cap B|}{|B|} = \frac{3}{4}.$$

■

This is a legitimate way to view the problem and it leads to the correct solution. However, this method has one very serious drawback – it requires us to change both our sample space and our probability function in order to carry out the computation. It would be preferable to have a method that allows us to work within the original framework of the sample space S and to talk about the “conditional probability” of A given that the result of the experiment will be an outcome in B . This is denoted as $P(A|B)$ and is read as “the (conditional) probability of A given B .”

Suppose S is a finite set of equally likely outcomes from a given experiment. Then for any two non-empty events A and B , the conditional probability of A given B is given by

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

This leads us to a formal definition of conditional probability for general sample spaces.

DEFINITION 1.3.2. (Conditional Probability) *Let S be a sample space with probability P . Let A and B be two events with $P(B) > 0$. Then the conditional probability of A given B written as $P(A|B)$ and is defined by*

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

This definition makes it possible to compute a conditional probability in terms of the original (unconditioned) probability function.

EXAMPLE 1.3.3. A pair of dice are thrown. If it is known that one die shows a 4, what is the probability the other die shows a 6?

Let B be the event that one of the dice shows a 4. So

$$B = \{(4, 1), (4, 2), (4, 3), (4, 4), (4, 5), (4, 6), (1, 4), (2, 4), (3, 4), (5, 4), (6, 4)\}.$$

Let A be the event that one of the dice is a 6. So

$$A = \{(6, 1), (6, 2), (6, 3), (6, 4), (6, 5), (6, 6), (1, 6), (2, 6), (3, 6), (4, 6), (5, 6)\}.$$

then

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

Hence

$$\begin{aligned} P(\text{one die shows 6} & \mid \text{one die shows 4}) \\ &= P(A|B) = \frac{P(A \cap B)}{P(B)} \\ &= \frac{P(\text{one die shows 6 and the other shows 4})}{P(\text{one die shows 4})} \\ &= \frac{2/36}{11/36} = \frac{2}{11}. \end{aligned}$$

■

In many applications, the conditional probabilities are implicitly defined within the context of the problem. In such cases, it is useful to have a method for computing non-conditional probabilities from the given conditional ones. Two such methods are given by the next results and the subsequent examples.

EXAMPLE 1.3.4. An economic model predicts that if interest rates rise, then there is a 60% chance that unemployment will increase, but that if interest rates do not rise, then there is only a 30% chance that unemployment will increase. If the economist believes there is a 40% chance that interest rates will rise, what should she calculate is the probability that unemployment will increase?

Let B be the event that interest rates rise and A be the event that unemployment increases. We know the values

$$P(B) = 0.4, P(B^c) = 0.6, P(A|B) = 0.6, \text{ and } P(A|B^c) = 0.3.$$

Using the axioms of probability and definition of conditional probability we have

$$\begin{aligned} P(A) &= P((A \cap B) \cup (A \cap B^c)) \\ &= P((A \cap B)) + P(A \cap B^c)) \\ &= P(A|B)P(B) + P(A|B^c)P(B^c) \\ &= 0.6(0.4) + 0.3(0.6) = 0.42. \end{aligned}$$

So there is a 42% chance that unemployment will increase. ■

THEOREM 1.3.5. Let A be an event and let $\{B_i : 1 \leq i \leq n\}$ be a disjoint collection of events for which $P(B_i) > 0$ for all i and such that $A \subset \bigcup_{i=1}^n B_i$. Suppose $P(B_i)$ and $P(A|B_i)$ are known. Then $P(A)$ may be computed as

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

Proof - The events $(A \cap B_i)$ and $(A \cap B_j)$ are disjoint if $i \neq j$ and

$$\bigcup_{i=1}^n (A \cap B_i) = A \cap (\bigcup_{i=1}^n B_i) = A.$$

So,

$$\begin{aligned} P(A) &= P\left(\bigcup_{i=1}^n (A \cap B_i)\right) \\ &= \sum_{i=1}^n P(A \cap B_i) \\ &= \sum_{i=1}^n P(A|B_i)P(B_i). \end{aligned}$$

■

A nearly identical proof holds when there are only countably many B_i (see Exercise 1.3.11).

EXAMPLE 1.3.6. Suppose we have coloured balls distributed in three boxes in quantities as given by the table below:

	Box 1	Box 2	Box 3
Red	4	3	3
Green	3	3	4
Blue	5	2	3

A box is selected at random. From that box a ball is selected at random. How likely is it that a red ball is drawn?

Let B_1 , B_2 , and B_3 be the events that Box 1, 2, or 3 is selected, respectively. Note that these events are disjoint and cover all possibilities in the sample space. Let R be the event that the selected ball is red. Then by Theorem 1.3.5,

$$\begin{aligned} P(R) &= P(R|B_1)P(B_1) + P(R|B_2)P(B_2) + P(R|B_3)P(B_3) \\ &= \frac{4}{12} \cdot \frac{1}{3} + \frac{3}{8} \cdot \frac{1}{3} + \frac{3}{10} \cdot \frac{1}{3} \\ &= \frac{121}{360}. \end{aligned}$$

■

EXAMPLE 1.3.7. (Polya's Urn Scheme) Suppose there is an urn that contains r red balls and b black balls. A ball is drawn at random and its colour noted. It is replaced with $c > 0$ balls of the same colour. The procedure is then repeated. For $j = 1, 2, \dots$, let R_j and B_j be the events that

the j -th ball drawn is red and black respectively. Clearly $P(R_1) = \frac{r}{b+r}$ and $P(B_1) = \frac{b}{b+r}$. When the first ball is replaced, c new balls will be added to the urn, so that when the second ball is drawn there will be $r+b+c$ balls available. From this it can easily be checked that $P(R_2|R_1) = \frac{r+c}{b+r+c}$ and $P(R_2|B_1) = \frac{r}{b+r+c}$. Noting that R_1 and B_1 are disjoint and together represent the entire sample space, $P(R_2)$ can be computed as

$$\begin{aligned} P(R_2) &= P(R_1)P(R_2|R_1) + P(B_1)P(R_2|B_1) \\ &= \frac{r}{b+r} \cdot \frac{r+c}{b+r+c} + \frac{b}{b+r} \cdot \frac{r}{b+r+c} \\ &= \frac{r(r+b+c)}{(r+b+c)(b+r)} \\ &= \frac{r}{b+r} = P(R_1). \end{aligned}$$

One can show that $P(R_j) = \frac{r}{b+r}$ for all $j \geq 1$.

The urn schemes were originally developed by George Polya (1887-1985). Various modifications to Polya's urn scheme are discussed in the exercises. ■

Above we have described how conditioning on an event B may be viewed as modifying the original probability based on the additional information provided by knowing that B has occurred. Frequently in applications we gain information more than once in the process of an experiment. The following theorem shows how to deal with such a situation.

THEOREM 1.3.8. *For an integer $n \geq 2$, let A_1, A_2, \dots, A_n be a collection of events for which $\bigcap_{j=1}^{n-1} A_j$ has positive probability. Then,*

$$P\left(\bigcap_{j=1}^n A_j\right) = P(A_1) \cdot \prod_{j=2}^n P(A_j | \bigcap_{k=1}^{j-1} A_k).$$

The proof of this theorem is left as Exercise 1.3.14, but we will provide a framework in which to make sense of the equality. Usually the events A_1, \dots, A_n are viewed as a sequence in time for which we know the probability of a given event provided that all of the others before it have already occurred. Then we can calculate $P(A_1 \cap A_2 \cap \dots \cap A_n)$ by taking the product of the values $P(A_1), P(A_2|A_1), P(A_3|A_1 \cap A_2), \dots, P(A_n|A_1 \cap \dots \cap A_{n-1})$.

EXAMPLE 1.3.9. A probability class has fifteen students - four seniors, eight juniors, and three sophomores. Three different students are selected at random to present homework problems. What is the probability the selection will be a junior, a sophomore, and a junior again, in that order?

Let A_1 be the event that the first selection is a junior. Let A_2 be the event that the second selection is a sophomore, and let A_3 be the event that the third selection is a junior. The problem asks for $P(A_1 \cap A_2 \cap A_3)$ which we can calculate using Theorem 1.3.8.

$$\begin{aligned} P(A_1 \cap A_2 \cap A_3) &= P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2) \\ &= \frac{8}{15} \cdot \frac{3}{14} \cdot \frac{7}{13} \\ &= \frac{4}{65}. \end{aligned}$$

■

1.3.1 Bayes' Theorem

It is often the case that we know the conditional probability of A given B , but want to know the conditional probability of B given A instead. It is possible to calculate one quantity from the other using a formula known as Bayes' theorem. We introduce this with a motivating example.

EXAMPLE 1.3.10. We return to Example 1.3.6. In that example we had three boxes containing balls given by the table below:

	Box 1	Box 2	Box 3
Red	4	3	3
Green	3	3	4
Blue	5	2	3

A box is selected at random. From the box a ball is selected at random. When we looked at conditional probabilities we saw how to determine the probability of an event such as $\{\text{the ball drawn is red}\}$. Now suppose we know the ball is red and want to determine the probability of the event $\{\text{the ball was drawn from box 3}\}$. That is, if R is the event that a red ball is chosen and if B_1 , B_2 , and B_3 are the events that boxes 1, 2, and 3 are selected, we want to determine the conditional probability $P(B_3|R)$. The difficulty is that while the conditional probabilities $P(R|B_1)$, $P(R|B_2)$, and $P(R|B_3)$ are easy to determine, calculating the conditional probability with the order of the events reversed is not immediately obvious.

Using the definition of conditional probability we have that

$$P(B_3|R) = \frac{P(B_3 \cap R)}{P(R)}.$$

We can rewrite

$$P(B_3 \cap R) = P(R|B_3)P(B_3) = \left(\frac{3}{10}\right)\left(\frac{1}{3}\right) = 0.1.$$

On the other hand, we can decompose the event R over which box was chosen. This is exactly what we did to solve Example 1.3.6 where we found that $P(R) = \frac{121}{360}$. Hence,

$$P(B_3|R) = \frac{P(B_3 \cap R)}{P(R)} = \frac{0.1}{121/360} = \frac{36}{121} \approx 0.298.$$

So if we know that a red ball was drawn, there is slightly less than a 30% chance that it came from Box 3. ■

In the above example the description of the experiment allowed us to determine $P(B_1)$, $P(B_2)$, $P(B_3)$, $P(R|B_1)$, $P(R|B_2)$, and $P(R|B_3)$. We were then able to use the definition of conditional probability to find $P(B_3|R)$. Such a computation can be done in general.

THEOREM 1.3.11. (Bayes' Theorem) Suppose A is an event, $\{B_i : 1 \leq i \leq n\}$ are a collection of disjoint events whose union contains all of A . Further assume that $P(A) > 0$ and $P(B_i) > 0$ for all $1 \leq i \leq n$. Then for any $1 \leq i \leq n$,

$$P(B_i|A) = \frac{P(A|B_i)P(B_i)}{\sum_{j=1}^n P(A|B_j)P(B_j)}.$$

Proof -

$$\begin{aligned}
 P(B_i | A) &= \frac{P(B_i \cap A)}{P(A)} = \frac{P(A | B_i)P(B_i)}{P(\bigcup_{j=1}^n A \cap B_j)} \\
 &= \frac{P(A | B_i)P(B_i)}{\sum_{j=1}^n P(A \cap B_j)} \\
 &= \frac{P(A | B_i)P(B_i)}{\sum_{j=1}^n P(A | B_j)P(B_j)}. \tag{1.3.1}
 \end{aligned}$$

■

The equation (1.3.1) is sometimes referred to as “Bayes’ formula” or “Bayes’ rule” as well. This is originally due to Thomas Bayes (1701-1761).

EXAMPLE 1.3.12. Shyam is randomly selected from the citizens of Hyderabad by the Health authorities. A laboratory test on his blood sample tells Shyam that he has tested positive for Swine Flu. It is found that 95% of people with Swine Flu test positive but 2% of people without the disease will also test positive. Suppose that 1% of the population has the disease. What is the probability that Shyam indeed has the Swine Flu ?

■

Consider the events $A = \{ \text{Shyam has Swine Flu} \}$ and $B = \{ \text{Shyam tested positive for Swine Flu} \}$. We are given:

$$P(B|A) = 0.95, P(B|A^c) = 0.02, \text{ and } P(A) = 0.01.$$

Using Bayes’ Theorem we have,

$$\begin{aligned}
 P(A|B) &= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^c)P(A^c)} \\
 &= \frac{(0.95)(0.01)}{(0.95)(0.01) + (0.02)(0.99)} \\
 &= 0.324
 \end{aligned}$$

Despite testing positive, there is only a 32.4 percent chance that Shyam has the disease.

EXERCISES

Ex. 1.3.1. There are two dice, one red and one blue, sitting on a table. The red die is a standard die with six sides while the blue die is tetrahedral with four sides, so the outcomes 1, 2, 3, and 4 are all equally likely. A fair coin is flipped. If that coin comes up heads, the red die will be rolled, but if the coin comes up tails the blue die will be rolled.

- (a) Find the probability that the rolled die will show a 1.
- (b) Find the probability that the rolled die will show a 6.

Ex. 1.3.2. A pair of dice are thrown. It is given that the outcome on one die is a 3. what is the probability that the sum of the outcomes on both dice is greater than 7?

Ex. 1.3.3. Box A contains four white balls and three black balls and Box B contains three white balls and five black balls.

- (a) Suppose a box is selected at random and then a ball is chosen from the box. If the ball drawn is black then what is the probability that it was from Box A?

- (b) Suppose instead that one ball is drawn at random from Box A and placed (unseen) in Box B. What is the probability that a ball now drawn from Box B is black?

Ex. 1.3.4. Tomorrow the weather will either be sunny, cloudy, or rainy. There is a 60% chance tomorrow will be cloudy, a 30% chance tomorrow will be sunny, and a 10% chance that tomorrow will be rainy. If it rains, I will not go on a walk. But if it is cloudy, there is a 90% chance I will take a walk and if it's sunny there is a 70% chance I will take a walk. If I take a walk on a cloudy day, there is an 80% chance I will walk further than five kilometers, but if I walk on a sunny day, there's only a 50% chance I will walk further than five kilometers. Using the percentages as given probabilities, answer the following questions:

- (a) How likely is it that tomorrow will be cloudy and I will walk over five kilometers?
- (b) How likely is it I will take a walk over five kilometers tomorrow?

Ex. 1.3.5. A box contains B black balls and W balls, where $W \geq 3, B \geq 3$. A sample of three balls is drawn at random with each drawn ball being discarded (not put back into the box) after it is drawn. For $j = 1, 2, 3$ let A_j denote the event that the ball drawn on the j^{th} draw is white. Find $P(A_1)$, $P(A_2)$ and $P(A_3)$.

Ex. 1.3.6. There are two sets of cards, one red and one blue. The red set has four cards - one that reads 1, two that read 2, and one that reads 3. An experiment involves flipping a fair coin. If the coin comes up heads a card will be randomly selected from the red set (and its number recorded) while if the coin comes up tails a card will be randomly selected from the blue set (and its number recorded). You can construct the blue set of cards in any way you see fit using any number of cards reading 1, 2, or 3. Explain how to build the blue set of cards to make each of the experimental outcomes 1, 2, 3 equally likely.

Ex. 1.3.7. There are three tables, each with two drawers. Table 1 has a red ball in each drawer. Table 2 has a blue ball in each drawer. Table 3 has a red ball in one drawer and a blue ball in the other. A table is chosen at random, then a drawer is chosen at random from that table. Find the conditional probability that Table 1 is chosen, given that a red ball is drawn.

Ex. 1.3.8. In the G.R.E advanced mathematics exam, each multiple choice question has 4 choices for an answer. A prospective graduate student taking the test knows the correct answer with probability $\frac{3}{4}$. If the student does not know the answer, she guesses randomly. Given that a question was answered correctly, find the conditional probability that the student knew the answer.

Ex. 1.3.9. You first roll a fair die, then toss as many fair coins as the number that showed on the die. Given that 5 heads are obtained, what is the probability that the die showed 5 ?

Ex. 1.3.10. Manish is a student in a probability class. He gets a note saying, "I've organized a probability study group tonight at 7pm in the coffee shop. Come if you want." The note is signed "Hannah". However, Manish has class with two different Hannahs and he isn't sure which one sent the note. He figures that there is a 75% chance that Hannah A. would have organized such a study group, but only a 25% chance that Hannah B. would have done so. However, he also figures that if Hannah A. had organized the group, there is an 80% chance that she would have planned to meet on campus and only a 20% chance that she would have planned to meet in the coffee shop. While if Hannah B. had organized the group there is a 10% chance she would have planned for it on campus and a 90% chance she would have chosen the coffee shop. Given all this information, determine whether it is more likely that Manish should think the note came from Hannah A. or from Hannah B.

Ex. 1.3.11. State and prove a version of

- (a) Theorem 1.3.5 when $\{B_i\}$ is a countably infinite collection of disjoint events.
- (b) Theorem 1.3.11 when $\{B_i\}$ is a countably infinite collection of disjoint events.

Ex. 1.3.12. A bag contains 100 coins. Sixty of the coins are fair. The rest are biased to land heads with probability p (where $0 \leq p \leq 1$). A coin is drawn at random from the bag and tossed.

- (a) Given that the outcome was a head what is the conditional probability that it is a biased coin?
- (b) Evaluate your answer to (a) when $p = 0$. Can you explain why this answer should be intuitively obvious?
- (c) Evaluate your answer to (a) when $p = \frac{1}{2}$. Can you explain why this answer should be fairly intuitive as well?
- (d) View your answer to part (a) as a function $f(p)$. Show that $f(p)$ is an increasing function when $0 \leq p \leq 1$. Give an interpretation of this fact in the context of the problem.

Ex. 1.3.13. An urn contains b black balls and r red balls. A ball is drawn at random. The ball is replaced into the urn along with c balls of its colour and d balls of the opposite colour. Then another random ball is drawn and the procedure is repeated.

- (a) What is the probability that the second ball drawn is a red ball?
- (b) Assume $c = d$. What is the probability that the second ball drawn is a black ball?
- (c) Still assuming $c = d$, what is the probability that the n^{th} ball drawn is a black ball?
- (d) Assume $c > 0$ and $d = 0$, what is the probability that the n^{th} ball drawn is a black ball?
- (e) Can you comment on the answers to (b) and/or (c) if the assumption that $c = d$ was removed?

Ex. 1.3.14. Use the following steps to prove Theorem 1.3.8.

- (a) Prove Theorem 1.3.8 for the $n = 2$ case. (Hint: The proof should follow immediately from the definition of conditional probability).
- (b) Prove Theorem 1.3.8 for the $n = 3$ case. (Hint: Rewrite the conditional probabilities in terms of ordinary probabilities).
- (c) Prove Theorem 1.3.8 generally. (Hint: One method is to use induction, and parts (a) and (b) have already provided a starting point).

1.4 INDEPENDENCE

In the previous section we have seen instances where the probability of an event may change given the occurrence of a related event. However it is instructive and useful to study the case of two events where the occurrence of one has no effect on the probability of the other. Such events are said to be “independent”.

EXAMPLE 1.4.1. Suppose we toss a coin three times. Then the sample space

$$S = \{hhh, hht, hth, htt, thh, tth, tht, ttt\}.$$

Define $A = \{hhh, hht, hth, htt\} = \{\text{the first toss is a head}\}$ and similarly define $B = \{hh, hht, thh, tht\} = \{\text{the second toss is a head}\}$. Note that $P(A) = \frac{1}{2} = P(B)$, while

$$P(A|B) = \frac{P(A \cap B)}{P(A)} = \frac{|A \cap B|}{|B|} = \frac{2}{4} = \frac{1}{2}$$

and

$$P(A|B^c) = \frac{P(A \cap B^c)}{P(B^c)} = \frac{|A \cap B^c|}{|B^c|} = \frac{2}{4} = \frac{1}{2}.$$

We have shown that $P(A) = P(A|B) = P(A|B^c)$. Therefore we conclude that the occurrence (or non-occurrence) of B has no effect on the probability of A . ■

This is the sort of condition we would want in a definition of independence. However, since defining $P(A|B)$ requires that $P(B) > 0$, our formal definition of “independence” will appear slightly different.

DEFINITION 1.4.2. (Independence) Two events A and B are independent if

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

EXAMPLE 1.4.3. Suppose we roll a die twice and denote as an ordered pair the result of the rolls. Suppose

$$E = \{\text{a six appears on the first roll}\} = \{(6, 1), (6, 2), (6, 3), (6, 4), (6, 5), (6, 6)\}$$

and

$$F = \{\text{a six appears on the second roll}\} = \{(1, 6), (2, 6), (3, 6), (4, 6), (5, 6), (6, 6)\}.$$

As $E \cap F = \{(6, 6)\}$, it is easy to see that

$$P(E \cap F) = \frac{1}{36}, P(E) = \frac{6}{36} = \frac{1}{6}, P(F) = \frac{6}{36} = \frac{1}{6}.$$

So E, F are independent as $P(E \cap F) = P(E)P(F)$. ■

Using the definition of conditional probability it is not hard to show (see Exercise 1.4.9) that if A and B are independent, and if $0 < P(B) < 1$ then

$$P(A|B) = P(A) = P(A|B^c). \quad (1.4.1)$$

If $P(A) > 0$ then the equations of (1.4.1) also hold with the roles of A and B reversed. Thus, independence implies four conditional probability equalities.

If we want to extend our definition of independence to three events A_1, A_2 , and A_3 , we would certainly want

$$P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2)P(A_3) \quad (1.4.2)$$

to hold. We would also want any pair of the three events to be independent of each other. It is tempting to hope that pairwise independence is enough to imply (1.4.2). However, consider the following example.

EXAMPLE 1.4.4. Suppose we toss a fair coin two times. Consider the three events $A_1 = \{hh, tt\}$, $A_2 = \{hh, ht\}$, and $A_3 = \{hh, th\}$. Then it is easy to calculate that

$$P(A_1) = P(A_2) = P(A_3) = \frac{1}{2},$$

$$P(A_1 \cap A_2) = P(A_1 \cap A_3) = P(A_2 \cap A_3) = \frac{1}{4}, \text{ and}$$

$$P(A_1 \cap A_2 \cap A_3) = \frac{1}{4}.$$

So even though A_1, A_2 and A_3 are pairwise independent this does not imply that they satisfy (1.4.2). ■

It may also be tempting to hope that (1.4.2) is enough to imply pairwise independence, but that is not true either (see Exercise 1.4.6). The root of the problem is that, unlike the two event case, (1.4.2) does not imply that equality holds if any of the A_i are replaced by their complements. One solution is to insist that the multiplicative equality hold for any intersection of the events or their complements, which gives us the following definition.

DEFINITION 1.4.5. (Mutual Independence) A finite collection of events A_1, A_2, \dots, A_n is mutually independent if

$$P(E_1 \cap E_2 \cap \dots \cap E_n) = P(E_1)P(E_2) \dots P(E_n). \quad (1.4.3)$$

whenever E_j is either A_j or A_j^c .

An arbitrary collection of events A_t where $t \in I$ for some index set I is mutually independent if every finite subcollection is mutually independent.

Thus, mutual independence of n events is defined in terms of 2^n equations. It is a fact (see Exercise 1.4.10) that if a collection of events is mutually independent, then so is any subcollection.

EXERCISES

Ex. 1.4.1. In the first semifinal of an international volleyball tournament Brazil has a 60% chance to beat Pakistan. In the other semifinal Poland has a 70% chance to beat Mexico. If the results of the two matches are independent, what is the probability that Pakistan will meet Poland in the tournament final?

Ex. 1.4.2. A manufacturer produces nuts and markets them as having 50mm radius. The machines that produce the nuts are not perfect. From repeated testing, it was established that 15% of the nuts have radius below 49mm and 12% have radius above 51mm. If two nuts are randomly (and independently) selected, find the probabilities of the following events:

- (a) The radii of both the nuts are between 49mm and 51mm;
- (b) The radius of at least one nut exceeds 51mm.

Ex. 1.4.3. Four tennis players (Avinash, Ben, Carlos, and David) play a single-elimination tournament with Avinash playing David and Ben playing Carlos in the first round and the winner of each of those contests playing each other in the tournament final. Below is the chart giving the percentage chance that one player will beat the other if they play. For instance, Avinash has a 30% chance of beating Ben if they happen to play.

	Avinash	Ben	Carlos	David
Avinash	-	30%	55%	40%
Ben	-	-	80%	45%
Carlos	-	-	-	15%
David	-	-	-	-

Suppose the outcomes of the games are independent. For each of the four players, determine the probability that player wins the tournament. Verify that the calculated probabilities sum to 1.

Ex. 1.4.4. Let A and B be events with $P(A) = 0.8$ and $P(B) = 0.7$.

- (a) What is the largest possible value of $P(A \cap B)$?
- (b) What is the smallest possible value of $P(A \cap B)$?
- (c) What is the value of $P(A \cap B)$ if A and B are independent?

Ex. 1.4.5. Suppose we toss two fair dice. Let E_1 denote the event that the sum of the dice is six. E_2 denote the event that sum of the dice equals seven. Let F denote the event that the first die equals four. Is E_1 independent of F ? Is E_2 independent of F ?

Ex. 1.4.6. Suppose a bowl has twenty-seven balls. One ball is black, two are white, and eight each are green, red, and blue. A single ball is drawn from the bowl and its color is recorded. Define

$$\begin{aligned} A &= \{\text{the ball is either black or green}\} \\ B &= \{\text{the ball is either black or red}\} \\ C &= \{\text{the ball is either black or blue}\} \end{aligned}$$

- (a) Calculate $P(A \cap B \cap C)$.
- (b) Calculate $P(A)P(B)P(C)$.
- (c) Are A , B , and C mutually independent? Why or why not?

Ex. 1.4.7. There are 150 students in the Probability 101 class. Of them, ninety are female, sixty use a pencil (instead of a pen), and thirty are wearing eye glasses. A student is chosen at random from the class. Define the following events:

$$\begin{aligned} A_1 &= \{\text{the student is a female}\} \\ A_2 &= \{\text{the student uses a pencil}\} \\ A_3 &= \{\text{the student is wearing eye glasses}\} \end{aligned}$$

- (a) Show that it is impossible for these events to be mutually independent.
- (b) Give an example to show that it may be possible for these events to be pairwise independent.

Ex. 1.4.8. When can an event be independent of itself? Do parts (a) and (b) below to answer this question.

- (a) Prove that if an event A is independent of itself then either $P(A) = 0$ or $P(A) = 1$.
- (b) Prove that if A is an event such that either $P(A) = 0$ or $P(A) = 1$ then A is independent of itself.

Ex. 1.4.9. This exercise explores the relationship between independence and conditional probability.

- (a) Suppose A and B are independent events with $0 < P(B) < 1$. Prove that $P(A|B) = P(A)$ and that $P(A|B^c) = P(A)$.
- (b) Suppose that A and B are independent events. Prove that A and B^c are also independent.
- (c) Suppose that A and B are events with $P(B) > 0$. Prove that if $P(A|B) = P(A)$, then A and B are independent.
- (d) Suppose that A and B are events with $0 < P(B) < 1$. Prove that if $P(A|B) = P(A)$, then $P(A|B^c) = P(A)$ as well.

Ex. 1.4.10. In this section we mentioned the following theorem: “If E_1, E_2, \dots, E_n is a collection of mutually independent events, then any subcollection of these events is mutually independent”. Follow the steps below to prove the theorem.

- (a) Suppose A, B , and C are mutually independent. In particular, this means that

$$P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C) \text{ and}$$

$$P(A \cap B \cap C^c) = P(A) \cdot P(B) \cdot P(C^c).$$

Use these two facts to conclude that A and B are pairwise independent.

- (b) Suppose E_1, E_2, \dots, E_n is a collection of mutually independent events. Prove that E_1, E_2, \dots, E_{n-1} is also mutually independent.
- (c) Use (b) and induction to prove the full theorem.

1.5 USING R FOR COMPUTATION

As we have already seen, and will see throughout this book, the general approach to solve problems in probability and statistics is to put them in an abstract mathematical framework. Many of these problems eventually simplify to computing some specific numbers. Usually these computations are simple and can be done using a calculator. For some computations however, a more powerful tool is needed. In this book, we will use a software called R to illustrate such computations. R is freely available open source software that runs on a variety of computer platforms, including Windows, Mac OS X, and GNU/Linux.

R is many different things to different people, but for our purposes, it is best to think of it as a very powerful calculator. Once you install and start R,¹ you will be presented with a prompt that looks like the “greater than” sign ($>$). You can type expressions that you want to evaluate here and press the Enter key to obtain the answer. For example,

```
> 9 / 44
[1] 0.2045455
> 0.6 * 0.4 + 0.3 * 0.6
[1] 0.42
> log(0.6 * 0.4 + 0.3 * 0.6)
[1] -0.8675006
```

¹Visit <https://www.r-project.org/> to download R and learn more about it.

It may seem odd to see a [1] at the beginning of each answer, but that is there for a good reason. R is designed for statistical computations, which often require working with a collection of numbers, which following standard mathematical terminology are referred to as *vectors*. For example, we may want to do some computations on a vector consisting of the first 5 positive integers. Specifically, suppose we want to compute the squares of these integers, and then sum them up. Using R, we can do

```
> c(1, 2, 3, 4, 5)^2
[1] 1 4 9 16 25
> sum(c(1, 2, 3, 4, 5)^2)
[1] 55
```

Here the construct `c(...)` is used to create a vector containing the first five integers. Of course, doing this manually is difficult for larger vectors, so another useful construct is `m:n` which creates a vector containing all integers from `m` to `n`. Just as we do in mathematics, it is also convenient to use symbols (called “variables”) to store intermediate values in long computations. For example, to do the same operations as above for the first 40 integers, we can do

```
> x <- 1:40
> x
[1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
[23] 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40
> x^2
[1] 1 4 9 16 25 36 49 64 81 100 121 144 169
[14] 196 225 256 289 324 361 400 441 484 529 576 625 676
[27] 729 784 841 900 961 1024 1089 1156 1225 1296 1369 1444 1521
[40] 1600
> sum(x^2)
[1] 22140
```

We can now guess the meaning of the number in square brackets at the beginning of each line in the output: when R prints a vector that spans multiple lines, it prefixes each line by the index of the first element printed in that line. The prefix appears for scalars too because R treats scalars as vectors of length one.

In the example above, we see two kinds of operations. The expression `x^2` is interpreted as an element-wise squaring operation, which means that the result will have the same length as the input. On the other hand, the expression `sum(x)` takes the elements of a vector `x` and computes their sum. The first kind of operation is called a *vectorized operation*, and most mathematical operations in R are of this kind.

To see how this can be useful, let us use R to compute factorials and binomial coefficients, which will turn up frequently in this book. Recall that the binomial coefficient

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

represents the number of ways of choosing k items out of n , where for any positive integer m , $m!$ is the product of the first m positive integers. Just as `sum(x)` computes the sum of the elements of `x`, `prod(x)` computes their product. So, we can compute $10!$ and $\binom{10}{4}$ as

```
> prod(1:10)
[1] 3628800

> prod(1:10) / (prod(1:4) * prod(1:6))

[1] 210
```

Unfortunately, factorials can quickly become quite big, and may be beyond R's ability to compute precisely even for moderately large numbers. For example, trying to compute $\binom{200}{4}$, we get

```
> prod(1:200)
[1] Inf

> prod(1:200) / (prod(1:4) * prod(1:196))

[1] NaN
```

The first computation yields `Inf` because at some point in the computation of the product, the result becomes larger than the largest number R can store (this is often called "overflowing"). The second computation essentially reduces to computing `Inf/Inf`, and the resulting `NaN` indicates that the answer is ambiguous. The trivial mathematical fact that

$$\log m! = \sum_{i=1}^m \log i$$

comes to our aid here because it lets us do our computations on much smaller numbers. Using this, we can compute

```
> logb <- sum(log(1:200)) - sum(log(1:4)) - sum(log(1:196))
> logb

[1] 17.98504

> exp(logb)

[1] 64684950
```

R actually has the ability to compute binomial coefficients built into it.

```
> choose(200, 4)

[1] 64684950
```

These named operations, such as `sum()`, `prod()`, `log()`, `exp()`, and `choose()`, are known as *functions* in R. They are analogous to mathematical functions in the sense that they map some inputs to an output. Vectorized functions map vectors to vectors, whereas summary functions like `sum()` and `prod()` map vectors to scalars. It is common practice in R to make functions vectorized whenever possible. For example, the `choose()` function is also vectorized:

```
> choose(10, 0:10)

[1] 1 10 45 120 210 252 210 120 45 10 1

> choose(10:20, 4)

[1] 210 330 495 715 1001 1365 1820 2380 3060 3876 4845

> choose(2:15, 0:13)
```

```
[1] 1 3 6 10 15 21 28 36 45 55 66 78 91 105
```

A detailed exposition of R is beyond the scope of this book. In this book, we will only use relatively basic R functions, which we will introduce as and when needed. There are many excellent introductions available for the interested reader, and the website accompanying this book (TODO) also contains some supplementary material. In particular, R is very useful for producing statistical plots, and most figures in this book are produced using R. We do not describe how to create these figures in the book itself, but R code to reproduce them is available on the website.

EXERCISES

Ex. 1.5.1. In R suppose we type in the following

```
> x <- c(-15, -11, -4, 0, 7, 9, 16, 23)
```

Find out the output of the built-in functions given below:

```
sum(x) length(x) mean(x) var(x) sd(x) max(x) min(x) median(x)
```

Ex. 1.5.2. Obtain a six-sided die, and throw it ten times, keeping a record of the face that comes up each time. Store these values in a vector variable x . Find the output of the built-in functions given in the previous exercise when applied to this vector.

Ex. 1.5.3. Use R to verify the calculations done in Example 1.2.4.

Ex. 1.5.4. We return to the Birthday Problem given in Exercise 1.2.12. Using R, calculate the Probability that at least two from a group of N people share the same birthday, for $N = 10, 12, 17, 26, 34, 40, 41, 45, 75, 105$.

2

SAMPLING AND REPEATED TRIALS

Consider an experiment and an event A within the sample space. We say the experiment is a success if an outcome from A occurs and failure otherwise. Let us consider the following examples:

Experiment	Sample Space	Event Description	Event A	P(A)
Toss a fair coin	$\{H, T\}$	Head appears	$\{H\}$	$\frac{1}{2}$
Roll a die	$\{1, 2, 3, 4, 5, 6\}$	Six appears	$\{6\}$	$\frac{1}{6}$
Roll a die	$\{1, 2, 3, 4, 5, 6\}$	A multiple of 3 appears	$\{3, 6\}$	$\frac{1}{3}$

In typical applications we would repeat an experiment several times independently and would be interested in the total number of successes achieved, a process that may be viewed as sampling from a large population. For instance, a manager in a factory making nuts and bolts, may devise an experiment to choose uniformly from a collection of manufactured bolts and call the experiment a success if the bolt is not defective. Then she would want to repeat such a selection every time and quantify the number of successes.

2.1 BERNOULLI TRIALS

We will now proceed to construct a mathematical framework for independent trials of an experiment where each trial is either a success or a failure. Let p be the probability of success at each trial. The sequence so obtained is called a sequence of Bernoulli trials with parameter p . The trials are named after James Bernoulli (1654-1705).

We will occasionally want to consider a single Bernoulli trial, so we will use the notation $Bernoulli(p)$ to indicate such a distribution. Since we are only interested in the result of the trial, we may view this as a probability on the sample space $S = \{\text{success}, \text{failure}\}$ where $P(\{\text{success}\}) = p$, but more often we will be interested in multiple independent trials. We discuss this in the next example.

EXAMPLE 2.1.1. Suppose we roll a die twice and ask how likely it is that we observe exactly one 6 between the two rolls. In the previous chapter (See Example 1.4.3) we would have viewed the sample space S as thirty-six equally likely outcomes, each of which was an ordered pair of results of the rolls. But since we are only concerned with whether the die roll is a 6 (success) or not a 6

(failure) we could also view it as two Bernoulli($\frac{1}{6}$) trials. Using notation from Example 1.4.3, note that $P(\text{success on the first roll}) = P(E) = \frac{1}{6}$ and $P(\text{success on the second roll}) = P(F) = \frac{1}{6}$. So

$$\begin{aligned} & P(\{\text{success, success}\}) \\ &= P(E \cap F) \\ &\quad (\text{using independence}) \\ &= P(E)P(F) \\ &= P(\text{success on the first roll}) \cdot P(\text{success on the second roll}) \\ &= \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}. \end{aligned}$$

We could alternately view S as having only four elements - (success,success), (success,failure), (failure,success), and (failure,failure). The four outcomes are not equally likely, but the fact that the trials are independent allows us to easily compute the probability of each. Through similar computations,

$$\begin{aligned} P(\{(success, failure)\}) &= 5/36, \\ P(\{(failure, success)\}) &= 5/36, \\ \text{and } P(\{(failure, failure)\}) &= 25/36. \end{aligned}$$

To complete the problem, the event of rolling exactly one 6 among the two dice requires exactly one success and exactly one failure. From the list above, this can happen in either of two orders, so the probability of observing exactly one 6 is $\frac{5}{36} + \frac{5}{36} = \frac{10}{36}$. ■

For any two real numbers a, b and any integer $n \geq 1$, it is well known that

$$(a+b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}. \quad (2.1.1)$$

This is the binomial expansion due to Blaise Pascal(1623-1662). It turns out when a and b are positive numbers with $a + b = 1$, the terms in the right hand side above have a probabilistic interpretation. We illustrate it in the example below.

EXAMPLE 2.1.2. After performing n independent Bernoulli(p) trials we are typically interested in the following questions:

- (a) What is the probability of observing exactly k successes?
- (b) What is the most likely number of successes?
- (c) How many attempts must be made before the first success is observed?
- (d) On average how many successes will there be?

Ans (a) - **Binomial(n,p)**: If $n = 1$, then the answer is clear, namely $P(\{\text{one success}\}) = p$ and $P(\{\text{zero successes}\}) = 1 - p$. For, $n > 1$ let $\omega = (\omega_1, \omega_2, \dots, \omega_n)$ be an n -tuple of outcomes. So we may view the sample space S as the set of all ω where each ω_i is allowed to be either “success” or “failure”. Let A_i represent either the event {the i^{th} trial is a success} or {the i^{th} trial is a failure}. Then by independence

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = \prod_{i=1}^n P(A_i). \quad (2.1.2)$$

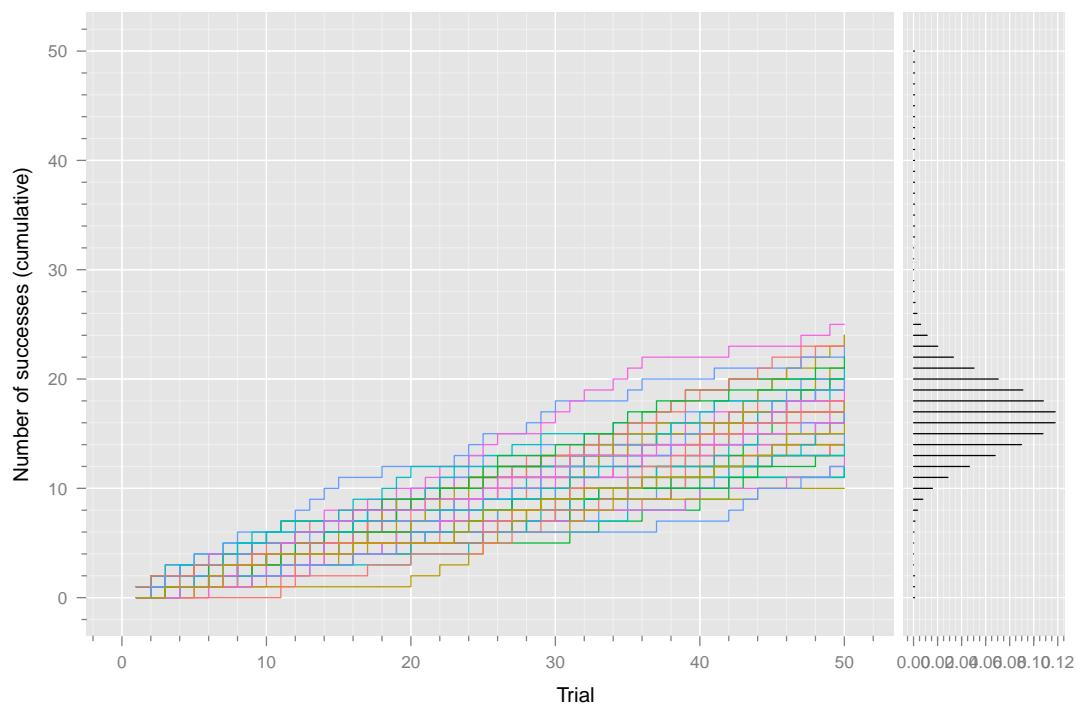


Figure 2.1: The Binomial distribution as number of successes in fifty Bernoulli ($\frac{1}{3}$) trials. The paths on the left count the cumulative successes in the fifty trials. The graph on the right show the actual probability given by the $\text{Binomial}(50, \frac{1}{3})$ distribution.

Let B_k denote the event that there are k successes among the n trials. Then

$$P(B_k) = \sum_{\omega \in B_k} P(\{\omega\}).$$

But if $\omega \in B_k$, then in notation (2.1.2), exactly k of the A_i represent success trials and the other $n - k$ represent the failure trials. The order in which the successes and failures appear does not matter since the probabilities are being multiplied together. So for every $\omega \in B_k$,

$$P(\{\omega\}) = p^k(1-p)^{n-k}.$$

Consequently, we have

$$P(B_k) = |B_k|p^k(1-p)^{n-k}.$$

But B_k is the event of all outcomes for which there are k successes and the number of ways in which k successes can occur in n trials is known to be $\binom{n}{k}$. Therefore, for $0 \leq k \leq n$,

$$P(B_k) = \binom{n}{k}p^k(1-p)^{n-k}. \quad (2.1.3)$$

Note that if we are only interested in questions involving the number of successes, we could ignore the set S described above and simply use $\{0, 1, 2, \dots, n\}$ as our sample space with $P(\{k\}) = \binom{n}{k}p^k(1-p)^{n-k}$. We call this a binomial distribution with parameters n and p (or a $\text{Binomial}(n, p)$ for short). It is also worth noting that the binomial expansion (2.1.1) shows

$$\sum_{k=0}^n \binom{n}{k}p^k(1-p)^{n-k} = (p + (1-p))^n = 1,$$

which simply provides additional confirmation that we have accounted for all possible outcomes in our list of Bernoulli trials. See Figure 2.1 for a simulated example of fifty replications of Bernoulli($\frac{1}{3}$) trials.

Ans (b) - Mode of a Binomial: The problem is trivial if $p = 0$ or $p = 1$, so assume $0 < p < 1$. Using the same notation for B_k as in part (a), pick a particular number of successes k for which $0 \leq k < n$. We want to determine the value of k that makes $P(B_k)$ as large as possible; such a value is called the “mode”. To find this value, it is instructive to compare the probability of $(k+1)$ successes to the probability of k successes –

$$\begin{aligned} \frac{P(B_{k+1})}{P(B_k)} &= \frac{\binom{n}{k+1}p^{k+1}(1-p)^{n-(k+1)}}{\binom{n}{k}p^k(1-p)^{n-k}} \\ &= \frac{n!}{(k+1)!(n-(k+1))!} \cdot \frac{k!(n-k)!}{n!} \cdot \frac{p^{k+1}(1-p)^{n-(k+1)}}{p^k(1-p)^{n-k}} \\ &= \frac{p}{1-p} \cdot \frac{n-k}{k+1}. \end{aligned}$$

If this ratio were to equal 1 we could conclude that $\{(k+1)\) successes\}$ was exactly as likely as $\{k\)$ successes $\}$. Similarly if the ratio were bigger than 1 we would know that $\{(k+1)\) successes\}$ was the more likely of the two and if the ratio were less than 1 we would see that $\{k\)$ successes $\}$ was the more likely case. Setting $\frac{P(B_{k+1})}{P(B_k)} \geq 1$ and solving for k yields the following sequence of equivalent inequalities:

$$\begin{aligned}
\frac{P(B_{k+1})}{P(B_k)} &\geq 1 \\
\frac{p}{1-p} \cdot \frac{n-k}{k+1} &\geq 1 \\
pn - pk &\geq k + 1 - pk - p \\
k &\leq p(n+1) - 1.
\end{aligned}$$

In other words if k starts at 0 and begins to increase, the probability of achieving exactly k successes will increase while $k < p(n+1) - 1$ and then will decrease once $k > p(n+1) - 1$. As a consequence the most likely number of successes is the integer value of k for which $k - 1 \leq p(n+1) - 1 < k$. This gives the critical value of $k = \lfloor p(n+1) \rfloor$, the greatest integer less than or equal to $p(n+1)$.

An unusual special case occurs if $p(n+1)$ is already an integer. Then the sequence of inequalities above is equality throughout, so if we let $k = \lfloor p(n+1) \rfloor = p(n+1)$ we find a ratio $P(B_k)/P(B_{k-1})$ exactly equal to 1. In this case $\{k-1 \text{ successes}\}$ and $\{k \text{ successes}\}$ share the distinction of being equally likely.

Ans (c) - Geometric(p): It is possible we could see the first success as early as the first trial and, in fact, the probability of this occurring is just p , the probability that the first trial is a success. The probability of the first success coming on the k^{th} trial requires that the first $k-1$ trials be failures and the k^{th} trial be a success. Let A_i be the event $\{\text{the } i^{\text{th}} \text{ trial is a success}\}$ and let C_k be the event $\{\text{the first success occurs on the } k^{\text{th}} \text{ trial}\}$. So,

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

As usual $P(A_i) = p$ and $P(A_i^c) = 1 - p$, so by independence

$$P(C_k) = P(A_1^c)P(A_2^c) \dots P(A_{k-1}^c)P(A_k) = (1-p)^{k-1}p$$

for $k > 0$. If we view these as probabilities of the outcomes of a sample space $\{1, 2, 3, \dots\}$, we call this a geometric distribution with parameter p (or a Geometric(p) for short).

Ans (d) - Average: This is a natural question to ask but it requires a precise definition of what we mean by “average” in the context of probability. We shall do this in Chapter 4 and return to answer (d) at that point in time. ■

Bernoulli trials may also be used to determine probabilities associated with who will win a contest that requires a certain number of individual victories. Below is an example applied to a “best two out of three” situation.

EXAMPLE 2.1.3. Jed and Sania play a tennis match. The match is won by the first player to win two sets. Sania is a bit better than Jed and she will win any given set with probability $\frac{2}{3}$. How likely is it that Sania will win the match? (Assume the results of each set are independent).

This can almost be viewed as three Bernoulli($\frac{2}{3}$) trials where we view a success as a set won by Sania. One problem with that perspective is that an outcome such as (win, win, loss) never occurs since two wins put an end to the match and the third set will never be played. Nevertheless, the same tools used to solve the earlier problem can be used for this one as well. Sania wins the match if she wins the first two sets (which happens with probability $\frac{4}{9}$). She also wins the match with either a (win, loss, win) or a (loss, win, win) sequence of sets, each of which has probability $\frac{4}{27}$ of occurring.

So the total probability of Sania winning the series is $\frac{4}{9} + \frac{4}{27} + \frac{4}{27} = \frac{20}{27}$.

Alternatively, it is possible to view this somewhat artificially as a genuine sequence of three Bernoulli($\frac{2}{3}$) trials where we pretend the players will play a third set even if the match is over by then. In effect the (win, win) scenario above is replaced by two different outcomes - (win, win, win) and (win, win, loss). Sania wins the match if she either wins all three sets (which has probability $\frac{8}{27}$) or if she wins exactly two of the three (which has probability $3 \cdot \frac{4}{27}$).

This perspective still leads us to the correct answer as $\frac{8}{27} + 3 \cdot \frac{4}{27} = \frac{20}{27}$. ■

2.1.1 Using R to compute probabilities

R can be used to compute probabilities of both the Binomial and Geometric distribution quite easily. We can compute them directly from the respective formulas. For example, with $n = 10$ and $p = 0.25$, all Binomial probabilities are given by

```
> k <- 0:5
> choose(5, k) * 0.25^k * 0.75^(5-k)
[1] 0.2373046875 0.3955078125 0.2636718750 0.0878906250 0.0146484375
[6] 0.0009765625
```

Similarly, the Geometric probabilities with $p = 0.25$ for $k = 0, 1, 2, \dots, 10$ are given by

```
> k <- 0:10
> 0.25 * 0.75^k
[1] 0.25000000 0.18750000 0.14062500 0.10546875 0.07910156 0.05932617
[7] 0.04449463 0.03337097 0.02502823 0.01877117 0.01407838
```

Actually, as both Binomial and Geometric are standard distributions, R has built-in functions to compute these probabilities as follows:

```
> dbinom(0:5, size = 5, prob = 0.25)
[1] 0.2373046875 0.3955078125 0.2636718750 0.0878906250 0.0146484375
[6] 0.0009765625
> dgeom(0:10, prob = 0.25)
[1] 0.25000000 0.18750000 0.14062500 0.10546875 0.07910156 0.05932617
[7] 0.04449463 0.03337097 0.02502823 0.01877117 0.01407838
```

EXERCISES

Ex. 2.1.1. Three dice are rolled. How likely is it that exactly one of the dice shows a 6?

Ex. 2.1.2. A fair die is rolled repeatedly.

- (a) What is the probability that the first 6 appears on the fifth roll?
- (b) What is the probability that no 6's appear in the first four rolls?
- (c) What is the probability that the second 6 appears on the fifth roll?

Ex. 2.1.3. Suppose that airplane engines operate independently in flight and fail with probability p ($0 \leq p \leq 1$). A plane makes a safe flight if at least half of its engines are running. Kingfisher Air lines has a four-engine plane and Paramount Airlines has a two-engine plane for a flight from Bangalore to Delhi. Which airline has the higher probability for a successful flight?

Ex. 2.1.4. Two intramural volleyball teams have eight players each. There is a 10% chance that any given player will not show up to a game, independently of any another. The game can be played if each team has at least six members show up. How likely is it the game can be played?

Ex. 2.1.5. Mark is a 70% free throw shooter. Assume each attempted free throw is independent of every other attempt. If he attempts ten free throws, answer the following questions.

- (a) How likely is it that Mark will make exactly seven of ten attempted free throws?
- (b) What is the most likely number of free throws Mark will make?
- (c) How do your answers to (a) and (b) change if Mark only attempts 9 free throws instead of 10?

Ex. 2.1.6. Continuing the previous exercise, Kalyani isn't as good a free throw shooter as Mark, but she can still make a shot 40% of the time. Mark and Kalyani play a game where the first one to sink a free throw is the winner. Since Kalyani isn't as skilled a player, she goes first to make it more fair.

- (a) How likely is it that Kalyani will win the game on her first shot?
- (b) How likely is it that Mark will win this game on his first shot? (Remember, for Mark even to get a chance to shoot, Kalyani must miss her first shot).
- (c) How likely is it that Kalyani will win the game on her second shot?
- (d) How likely is it that Kalyani will win the game?

Ex. 2.1.7. Recall from the text above that the R code

```
> dbinom(0:5, size = 5, prob = 0.25)
[1] 0.2373046875 0.3955078125 0.2636718750 0.0878906250 0.0146484375
[6] 0.0009765625
```

produces a vector of six outputs corresponding to the probabilities that a Binomial(5, 0.25) distribution takes on the six values 0-5. Specifically, the output indicates that the probability of the value 0 is approximately 0.2373046875, the probability of the value 1 is approximately 0.3955078125 and so on. In Example 2.1.2 we derived a formula for the most likely outcome of such a distribution. In the case of a Binomial(5, 0.25) that formula gives the result $\lfloor (5+1)(0.25) \rfloor = 1$. We could have verified this via the R output above as well, since the second number on the list is the largest of the probabilities.

- (a) Use the formula from example 2.1.2 to find the most likely outcome of a Binomial(7, 0.34) distribution.
- (b) Type an appropriate command into R to produce a vector of values corresponding to the probabilities that a Binomial(7, 0.34) distribution takes on the possible values in its range. Use this list to verify your answer to part (a).
- (c) Use the formula from Example 2.1.2 to find the most likely outcome of a Binomial(8, 0.34) distribution.
- (d) Type an appropriate command into R to produce a vector of values corresponding to the probabilities that a Binomial(8, 0.34) distribution takes on the possible values in its range. Use this list to verify your answer to part (c).

Ex. 2.1.8. It is estimated that 0.8% of a large shipment of eggs to a certain supermarket are cracked. The eggs are packaged in cartons, each with a dozen eggs, with the cracked eggs being randomly distributed. A restaurant owner buys 10 cartons from the supermarket. Call a carton "defective" if it contains at least one cracked egg.

- (a) If she notes the number of defective cartons, what are the possible outcomes for this experiment?
- (b) If she notes the total number of cracked eggs, what are the possible outcomes for this experiment?
- (c) How likely is it that she will find exactly one cracked egg among all of her cartons?
- (d) How likely is it that she will find exactly one defective carton?
- (e) Explain why your answer to (d) is close to, but slightly larger than, than your answer to (c).
- (f) What is the most likely number of cracked eggs she will find among her cartons?
- (g) What is the most likely number of defective cartons she will find?
- (h) How do you reconcile your answers to parts (g) and (h)?

Ex. 2.1.9. Steve and Siva enter a bar with \$30 each. A round of drinks cost \$10. For each round, they roll a die. If the roll is even, Steve pays for the round and if the roll is odd, Siva pays for it. This continues until one of them runs out of money.

- (a) What is the Probability that Siva runs out of money?
- (b) What is the Probability that Siva runs out of money if Steve has cheated by bringing a die that comes up even only 40% of the time?

Ex. 2.1.10. Let $0 < p < 1$. Show that the mode of a Geometric(p) distribution is 1.

Ex. 2.1.11. Scott is playing a game where he rolls a standard die until it shows a 6. The number of rolls needed therefore has a Geometric($\frac{1}{6}$) distribution. Use the appropriate R commands to do the following:

- (a) Produce a vector of values for $j = 1, \dots, 6$ corresponding to the probabilities that it will take Scott j rolls before he observes a 6.
- (b) Scott figures that since each roll has a $\frac{1}{6}$ probability of producing a 6, he's bound to get that result at some point after six rolls. Use the results from part (a) to determine the probability that Scott's expectations are met and a 6 will show up in one his first six rolls.

Ex. 2.1.12. Suppose a fair coin is tossed n times. Compute the following:

- (a) $P(\{4 \text{ heads occur}\} | \{3 \text{ or } 4 \text{ heads occur}\})$;
- (b) $P(\{k - 1 \text{ heads occur}\} | \{k - 1 \text{ or } k \text{ heads occur}\})$; and
- (c) $P(\{k \text{ heads occur}\} | \{k - 1 \text{ or } k \text{ heads occur}\})$.

Ex. 2.1.13. At a basketball tournament, each round is on a “best of seven games” basis. That is, Team I and Team 2 play until one of the teams has won four games. Suppose each game is won by Team I with probability p , independently of all previous games. Are the events $A = \{\text{Team I wins the round}\}$ and $B = \{\text{the round lasts exactly four games}\}$ independent?

Ex. 2.1.14. Two coins are sitting on a table. One is fair and the other is weighted so that it always comes up heads.

- (a) If one coin is selected at random (each equally likely) and flipped, what is the probability the result is heads?

- (b) One coin is selected at random (each equally likely) and flipped five times. Each flip shows heads. Given this information about the coin flip results, what is the conditional probability that the selected coin was the fair one?

Ex. 2.1.15. For $0 < p < 1$ we defined the geometric distribution as a probability on the set $\{1, 2, 3, \dots\}$ for which $P(\{k\}) = p(1-p)^{k-1}$. Show that these outcomes account for all possibilities by demonstrating that $\sum_{k=1}^{\infty} P(\{k\}) = 1$.

Ex. 2.1.16. The geometric distribution described the waiting time to observe a single success. A “negative binomial” distribution with parameters n and p ($\text{NegBinomial}(n, p)$) is defined the number of Bernoulli(p) trials needed before observing n successes. The following problem builds toward calculating some associated probabilities.

- (a) If a fair die is rolled repeatedly and a number is recorded equal to the number of rolls until the second 6 is observed, what is the sample space of possible outcomes for this experiment?
- (b) For k in the sample space you identified in part (a), what is $P(\{k\})$?
- (c) If a fair die is rolled repeatedly and a number is recorded equal to the number of rolls until the n^{th} 6 is observed, what is the sample space of possible outcomes for this experiment?
- (d) For k in the sample space you identified in part (c), what is $P(\{k\})$?
- (e) If a sequence of Bernoulli(p) trials (with $0 < p < 1$) is performed and a number is recorded equal to the number of trials until the n^{th} success is observed, what is the sample space of possible outcomes for this experiment?
- (f) For k in the sample space you identified in part (e), what is $P(\{k\})$?
- (g) Show that you have accounted for all possibilities in part (f) by showing

$$\sum_{k \in S} P(\{k\}) = 1.$$

2.2 POISSON APPROXIMATION

Calculating binomial probabilities can be challenging when n is large. Let us consider the following example:

EXAMPLE 2.2.1. A small college has 1460 students. Assume that birthrates are constant throughout the year and that each year has 365 days. What is the probability that five or more students were born on Independence day?

The probability that any given student was born on Independence day is $\frac{1}{365}$. So the exact probability is

$$1 - \sum_{k=0}^4 \binom{1460}{k} \left(\frac{1}{365}\right)^k \left(\frac{364}{365}\right)^{1460-k}.$$

Repeatedly dealing with large powers of fractions or large combinatorial computations is not so easy, so it would be convenient to find a faster way to estimate such a probability. ■

The example above can be thought of as a series of Bernoulli trials where a success means finding a student whose birthday is Independence day. In this case p is small ($\frac{1}{365}$) and n is large (1460). To approximate we will consider a limiting procedure where $p \rightarrow 0$ and $n \rightarrow \infty$, but with limits carried out in such a way that np is held constant. The computation below is called a Poisson approximation.

THEOREM 2.2.2. Let $\lambda > 0$, $k \geq 1$, $n \geq \lambda$ and $p = \frac{\lambda}{n}$. Defining A_k as

$$A_k = \{k \text{ successes in } n \text{ Bernoulli}(p) \text{ Trials}\},$$

it then follows that

$$\lim_{n \rightarrow \infty} P(A_k) = \frac{e^{-\lambda} \lambda^k}{k!}. \quad (2.2.1)$$

Proof -

$$\begin{aligned} P(A_k) &= \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \frac{n(n-1)\dots(n-k+1)}{k!} \frac{\lambda^k}{n^k} \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \frac{\lambda^k}{k!} \frac{n(n-1)\dots(n-k+1)}{n^k} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k} \\ &= \frac{\lambda^k}{k!} \left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{k-1}{n}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k} \\ &= \frac{\lambda^k}{k!} \prod_{r=1}^{k-1} \left(1 - \frac{r}{n}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k}. \end{aligned}$$

Standard limit results imply that

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(1 - \frac{r}{n}\right) &= 1 \quad \text{for all } r \geq 1; \\ \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-k} &= 1 \quad \text{for all } \lambda \geq 0, k \geq 1; \text{ and} \\ \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n &= e^{-\lambda} \quad \text{for all } \lambda \geq 0. \end{aligned}$$

As $P(A_k)$ is a finite product of such expressions, the result is now immediate using the properties of limits. ■

Returning to Example 2.2.1 and using the above approximation, we would take $\lambda = pn = \frac{1460}{365} = 4$. So if E is the event {five or more Independence day birthdays},

$$\begin{aligned} P(E) &= 1 - \sum_{k=0}^4 \binom{1460}{k} \left(\frac{1}{365}\right)^k \left(\frac{364}{365}\right)^{1460-k} \\ &\approx 1 - \left[e^{-4} + 4e^{-4} + \frac{4^2}{2}e^{-4} + \frac{4^3}{6}e^{-4} + \frac{4^4}{24}e^{-4} \right]. \end{aligned}$$

Calculation demonstrates this is a good approximation. To seven digits of accuracy, the correct value is 0.37116294 while the Poisson approximation gives an answer of 0.37116306. These can be obtained using R as follows:

```
> 1 - sum(dbinom(0:4, size = 1460, prob = 1/365))
[1] 0.3711629
> lambda <- 1460 / 365
> 1 - sum(exp(-lambda) * lambda^(0:4) / factorial(0:4))
[1] 0.3711631
```

It also turns out that the right hand side of (2.2.1) defines a probability on the sample space of non-negative integers. The distribution is named after Siméon Poisson (1781-1840).

Poisson (λ): Let $\lambda \geq 0$ and $S = \{0, 1, 2, 3, \dots\}$ with probability P given by

$$P(\{k\}) = \frac{e^{-\lambda} \lambda^k}{k!}$$

for $k \in S$. This distribution is called Poisson with parameter λ (or $\text{Poisson}(\lambda)$ for short).

As with Binomial and Geometric, R has a built-in function to evaluate Poisson probabilities as well. An alternative to the calculation above is:

```
> 1 - sum(dpois(0:4, lambda = 1460 / 365))
[1] 0.3711631
```

It is important to note that for this approximation to work well, p must be small and n must be large. For example, we may modify our question as follows:

EXAMPLE 2.2.3. A class has 48 students. Assume that birthrates are constant throughout the year and that each year has 365 days. What is the probability that five or more students were born in September? ■

The correct answer to this question is

```
> 1 - sum(dbinom(0:4, size = 48, prob = 1/12))
0.3710398
```

However, the Poisson approximation remains unchanged at 0.3711631, because $np = 48/12 = 1460/365 = 4$, and only matches the correct answer up to 3 digits rather than 6. Figure 2.2 shows a point-by-point approximation of both Binomial distributions by Poisson.

At this point we have defined many named distributions. Frequently a problem will require the use of more than one of these as evidenced in the next example.

EXAMPLE 2.2.4. A computer transmits three digital messages of 12 million bits of information each. Each bit has a probability of one one-billionth that it will be incorrectly received, independent of all other bits. What is the probability that at least two of the three messages will be received error free?

Since $n = 12,000,000$ is large and since $p = \frac{1}{1,000,000,000}$ is small it is appropriate to use a Poisson approximation where $\lambda = np = 0.012$. A message is error free if there isn't a single misread bit, so the probability that a given message will be received without an error is $e^{-0.012}$.

Now we can think of each message being like a Bernoulli trial with probability $e^{-0.012}$, so the number of messages correctly received is then like a $\text{Binomial}(3, e^{-0.012})$. Therefore the probability of receiving at least two error-free messages is

$$\binom{3}{0} (e^{-0.012})^0 (1 - e^{-0.012})^3 + \binom{3}{1} (e^{-0.012})^1 (1 - e^{-0.012})^2 \approx 0.9996.$$

There is about a 99.96% chance that at least two of the messages will be correctly received. ■

EXERCISES

Ex. 2.2.1. Do the problems below to familiarize yourself with the “sum” command in R.

- (a) If a fair coin is tossed 100 times, what is the probability exactly 55 of the tosses show heads?
- (b) Example 2.2.3 showed how to use R to add the probabilities of a range of outcomes for common distributions. Use this code as a guide to calculate the probability at least 55 tosses show heads.

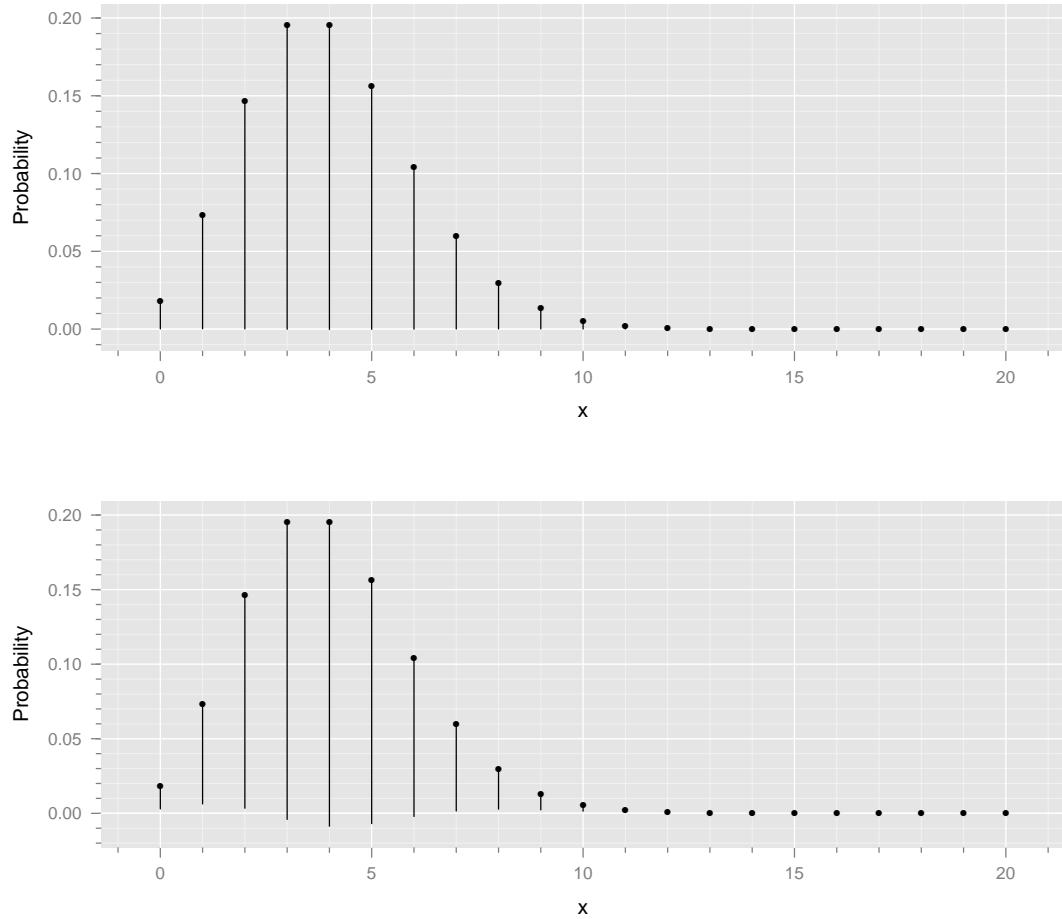


Figure 2.2: The Poisson approximation to the Binomial distribution. In both plots above, the points indicate Binomial probabilities for $k = 0, 1, 2, \dots, 20$; the top plot for $\text{Binomial}(1460, \frac{1}{365})$, and the bottom for $\text{Binomial}(48, \frac{1}{12})$. The lengths of the vertical lines, “hanging” from the points, represent the corresponding probabilities for $\text{Poisson}(4)$. For a good approximation, the bottom of the hanging lines should end up at the x-axis. As we can see, this happens in the top plot but not for the bottom plot, indicating that $\text{Poisson}(4)$ is a good approximation for the first Binomial distribution, but not as good for the second.

Ex. 2.2.2. Consider an experiment described by a Poisson($\frac{1}{2}$) distribution and answer the following questions.

- (a) What is the probability the experiment will produce a result of 0?
- (b) What is the probability the experiment will produce a result larger than 1?

Ex. 2.2.3. Suppose we perform 500 independent trials with probability of success being 0.02.

- (a) Use R to compute the probability that there are six or fewer successes. Obtain a decimal approximation accurate to five decimal places.
- (b) Use the Poisson approximation to estimate the probability that there are six or fewer successes and compare it to your answer to (a).

Now suppose we perform 5000 independent trials with probability of success being 0.002.

- (c) Use R to compute the probability that there are six or fewer successes. Obtain a decimal approximation accurate to five decimal places.
- (d) Use the Poisson approximation to estimate the probability that there are six or fewer successes and compare it to your answer to (c).
- (e) Which approximation (b) or (d) is more accurate? Why?

Ex. 2.2.4. For a certain daily lottery, the probability is $\frac{1}{10000}$ that you will win. Suppose you play this lottery every day for three years. Use the Poisson approximation to estimate the chance that you will win more than once.

Ex. 2.2.5. A book has 200 pages. The number of mistakes on each page has a Poisson(1) distribution, and is independent of the number of mistakes on all other pages.

- (a) What is the chance that there are at least 2 mistakes on the first page?
- (b) What is the chance that at least eight of the first ten pages are free of mistakes?

Ex. 2.2.6. Let $\lambda > 0$. For the problems below, assume the probability space is a Poisson(λ) distribution.

- (a) Let k be a non-negative integer. Calculate the ratio $\frac{P(\{k+1\})}{P(\{k\})}$.
- (b) Use (a) to calculate the mode of a Poisson(λ).

Ex. 2.2.7. A number is to be produced as follows. A fair coin is tossed. If the coin comes up heads the number will be the outcome of an experiment corresponding to a Poisson(1) distribution. If the coin comes up tails the number will be the outcome of an experiment corresponding to a Poisson(2) distribution. Given that the number produced was a 2, determine the conditional probability that the coin came up heads.

Ex. 2.2.8. Suppose that the number of earthquakes that occur in a year in California has a Poisson distribution with parameter λ . Suppose that the probability that any given earthquake has magnitude at least 6 on the Richter scale is p .

- (a) Given that there are exactly n earthquakes in a year, find an expression (in terms of n and p) for the conditional probability that exactly one of them is magnitude at least 6.

- (b) Find an expression (in terms of λ and p) for the probability that there will be exactly one earthquake of magnitude at least 6 in a year.
- (c) Find an expression (in terms of n , λ , and p) for the probability that there will be exactly k earthquakes of magnitude at least 6 in a year.

Ex. 2.2.9. We defined a Poisson distribution as a probability on $S = \{0, 1, 2, \dots\}$ for which

$$P(\{k\}) = \frac{e^{-\lambda} \lambda^k}{k!},$$

for $k \geq 1$. Prove that this completely accounts for all possibilities by proving that

$$\sum_{k=0}^{\infty} \frac{e^{-\lambda} \lambda^k}{k!} = 1.$$

(Hint: Consider the power series expansion of the exponential function).

Ex. 2.2.10. Consider n vertices labeled $\{1, 2, \dots, n\}$. Corresponding to each distinct pair $\{i, j\}$ we perform an independent Bernoulli (p) experiment and insert an edge between i and j with probability p . The graph constructed this way is denoted as $G(n, p)$.

- (a) Let $1 \leq i \leq n$. We say j is a neighbour of i if there is an edge between i and j . For some $1 \leq k \leq n$ determine the probability that i has k neighbours ?
- (b) Let $\lambda > 0$ and n large enough so that $0 < p = \frac{\lambda}{n} < 1$ and let $A_k = \{ \text{vertex } 1 \text{ has } k \text{ neighbours} \}$ what is the

$$\lim_{n \rightarrow \infty} P(A_k)?$$

2.3 SAMPLING WITH AND WITHOUT REPLACEMENT

Imagine a small town with 5000 residents, exactly 1000 of whom are under the age of eighteen. Suppose we randomly select four of these residents and ask how many of the four are under the age of eighteen. There is some ambiguity in how to interpret this idea of selecting four residents. One possibility is “sampling with replacement” where each selection could be any of the 5000 residents and the selections are all genuinely independent. With this interpretation, the sample is simply a series of four independent Bernoulli($\frac{1}{5}$) trials, in which case the answer may be found using techniques from the previous sections. Note, however, that the assumption of independence allows for the possibility that the same individual will be chosen two or more times in separate trials. This is a situation that might seem peculiar when we think about choosing four people from a population of 5000, since we may not have four different individuals at the end of the process. To eliminate this possibility consider “sampling without replacement” where it is assumed that if an individual is chosen for inclusion in the sample, that person is no longer available to be picked in a later selection. Equivalently we can consider all possible groups of four which might be selected and view each grouping as equally likely. This change means the problem can no longer be solved by viewing the situation as a series of independent Bernoulli trials. Nevertheless, other tools that have been previously developed will serve to answer this new problem.

EXAMPLE 2.3.1. For the town described above, what is the probability that, of four residents randomly selected (without replacement), exactly two of them will be under the age of eighteen?

Since we are selecting four residents from the town of 5000, there are $\binom{5000}{4}$ ways this may be done. If each of these is equally likely, the desired probability may be calculated by determining

how many of these selections result in exactly two people under the age of eighteen. This requires selecting two of the 1000 who are in that younger age group and also selecting two of the 4000 who are older. So there are $\binom{1000}{2} \binom{4000}{2}$ ways to make such choices and therefore the probability of selecting exactly two residents under age eighteen is $\binom{1000}{2} \binom{4000}{2} / \binom{5000}{4}$.

It is instructive to compare this to the solution if it is assumed the selection is done with replacement. In that case, the answer is the simply the probability that a Binomial(4, $\frac{1}{5}$) produces a result of two. From the previous sections, the answer is $\binom{4}{2} (\frac{1}{5})^2 (\frac{4}{5})^2$.

To compare these answers we give decimal approximations of both. To six digits of accuracy

$$\frac{\binom{1000}{2} \binom{4000}{2}}{\binom{5000}{4}} \approx 0.153592 \quad \text{and} \quad \binom{4}{2} \left(\frac{1}{5}\right)^2 \left(\frac{4}{5}\right)^2 = 0.1536,$$

so while the two answers are not equal, they are very close. This is a reflection of an important fact in statistical analysis – when samples are small relative to the size of the populations they came from, the two methods of sampling give very similar results. ■

2.3.1 The Hypergeometric Distribution

Analyzing such problems more generally, consider a population of N people. Suppose r of these N share a common characteristic and the remaining $N - r$ do not have this characteristic. We take a sample of size m (without replacement) from the population and count the number of people among the sample that have the specified characteristic. This experiment is described by probabilities known as a hypergeometric distribution. Notice that the largest possible result is $\min\{m, r\}$ since the number cannot be larger than the size of the sample nor can it be larger than the number of people in the population with the characteristic. On the other extreme, it may be that the sample is so large it is guaranteed to select some people with the characteristic simply because the number of people without has been exhausted. More precisely, for every selection over $N - r$ in the sample we are guaranteed to select at least one person who has the characteristic. So the minimum possible result is the larger of 0 or $(m - (N - r))$.

HyperGeo(N, r, m): Let r and m be non-negative integers and let N be an integer with $N > \max\{r, m\}$. Let S be the set of integers ranging from $\max\{0, m - (N - r)\}$ to $\min\{m, r\}$ inclusive with probability P given by

$$P(\{k\}) = \frac{\binom{r}{k} \binom{N-r}{m-k}}{\binom{N}{m}}$$

for $k \in S$. Such a distribution is called hypergeometric with parameters N , r , and m (or HyperGeo(N, r, m)).

Of course, R can be used to compute hypergeometric probabilities as well. Example 2.3.1 can be phrased in terms of a HyperGeo(5000, 1000, 4) distribution, with $P(\{k\})$ being the desired answer. This probability can be computed as:

```
> dhyper(2, 4000, 1000, 4)
[1] 0.1535923
```

Note however, that instead of N , the parameter used by R is $N - r$.

2.3.2 Hypergeometric Distributions as a Series of Dependent Trials

It is also possible (and useful) to view sampling without replacement as a series of dependent Bernoulli trials for which each trial reduces the possible outcomes of subsequent trials. In this case

each trial is described in terms of conditional probabilities based on the results of the preceding observations. We illustrate this by revisiting the previous example.

Example 2.3.1 Continued: We first solved this problem by considering every group of four as equally likely to be selected. Now consider the sampling procedure as a series of four separate Bernoulli trials where a success corresponds to the selection of a person under eighteen and a failure as the selection of someone older. We still want to determine the probability that a sample of size four will produce exactly two successes. One complication with this perspective is that the successes and failures could come in many different orders, so first consider the event where the series of selections follow the pattern “success-success-failure-failure”. More precisely, for $j = 1, 2, 3, 4$ let

$$A_j = \{\text{The } j^{\text{th}} \text{ selection is a person younger than eighteen}\}.$$

Clearly $P(A_1) = \frac{1000}{5000}$. Given that the first selection is someone under eighteen, there are now only 4999 people remaining to choose among, and only 999 of them are under eighteen. Therefore $P(A_2|A_1) = \frac{999}{4999}$. Continuing with that same reasoning,

$$P(A_3^c|A_1 \cap A_2) = \frac{4000}{4998}$$

and

$$P(A_4^c|A_1 \cap A_2 \cap A_3^c) = \frac{3999}{4997}.$$

From those values, Theorem 1.3.8 may be used to calculate

$$\begin{aligned} P(\text{success - success} - \text{failure - failure}) &= P(A_1 \cap A_2 \cap A_3^c \cap A_4^c) \\ &= \frac{1000}{5000} \cdot \frac{999}{4999} \cdot \frac{4000}{4998} \cdot \frac{3999}{4997}. \end{aligned}$$

Next we must account for the fact that this figure only considers the case where the two younger people were chosen as the first two selections. There are $\binom{4}{2}$ different orderings that result in two younger and two older people, and it happens that each of these has the same probability calculated above. For example,

$$\begin{aligned} P(\text{failure - success} - \text{success - failure}) &= P(A_1^c \cap A_2 \cap A_3 \cap A_4^c) \\ &= \frac{4000}{5000} \cdot \frac{1000}{4999} \cdot \frac{999}{4998} \cdot \frac{3999}{4997}. \end{aligned}$$

The individual fractions are different, but their product is the same. This will always happen for different orderings of a specific number of successes since the denominators (5000 through 4997) reflect the steady reduction of one available choice with each additional selection. Similarly the numerators (1000 and 999 together with 4000 and 3999) reflect the number of people available from each of the two different categories and their reduction as previous choices eliminate possible candidates. Therefore the total probability is the product of the number of orderings and the probability of each ordering.

$$P(\text{two under eighteen}) = \binom{4}{2} \cdot \frac{4000}{5000} \cdot \frac{1000}{4999} \cdot \frac{999}{4998} \cdot \frac{3999}{4997}.$$

We leave it to the reader to verify that this is equal to $\binom{1000}{2} \binom{4000}{2} / \binom{5000}{4}$, the answer we found when we originally solved the problem via a different method.

The following theorem generalizes this previous example.

THEOREM 2.3.2. Let S be a sample space with a hypergeometric distribution with parameters N , r , and m . Then $P(\{k\})$ equals

$$\binom{m}{k} \left[\frac{r}{N} \frac{r-1}{N-1} \cdots \frac{r-(k-1)}{N-(k-1)} \right] \left[\frac{N-r}{N-k} \frac{N-r-1}{N-k-1} \cdots \frac{N-r-(m-1-k)}{N-(m-1)} \right]$$

for any $k \in S$.

Proof- Following the previous example as a model, this can be proven by viewing the hypergeometric distribution as a series of dependent trials. The first k fractions are the probabilities the first k trials each result in successes conditioned on the successes of the preceding trials. The remaining $m - k$ fractions are the conditional probabilities the remaining trials result in failures. The leading factor of $\binom{m}{k}$ accounts for the number of different patterns of k successes and $m - k$ failures, each of which is equally likely. It is also possible to prove the equality directly using combinatorial identities and we leave this as Exercise 2.3.4. ■

2.3.3 Binomial Approximation to the Hypergeometric Distribution

We saw with Example 2.3.1 that sampling with and without replacement may give very similar results. The following theorem makes a precise statement to this effect.

THEOREM 2.3.3. Let N , m , and r be positive integers for which $m < r < N$ and let k be a positive integer between 0 and m . Define

$$p = \frac{r}{N}, \quad p_1 = \frac{r-k}{N-k}, \quad \text{and} \quad p_2 = \frac{r-k}{N-m}.$$

Letting H denote the probability that a hypergeometric distribution with parameters N , r , and m takes on the value k , the following inequalities give bounds on this probability:

$$\binom{m}{k} p_1^k (1-p_1)^{m-k} < H \leq \binom{m}{k} p^k (1-p)^{m-k}.$$

Proof- The inequalities may be verified by comparing p , p_1 , and p_2 to the fractions from Theorem 2.3.2. Specifically note that the k fractions

$$\frac{r}{N}, \frac{r-1}{N-1}, \dots, \frac{r-(k-1)}{N-(k-1)}$$

are all less than or equal to p . Likewise the $m - k$ fractions

$$\frac{N-r}{N-k}, \frac{N-r-1}{N-k-1}, \dots, \frac{N-r-(m-1-k)}{N-(m-1)}$$

are all less than or equal to $\frac{N-r}{N-k}$ which itself equals $1 - p_1$. Combining these facts proves the right hand inequality. The left hand inequality may be similarly shown by noting that the fractions

$$\frac{r}{N}, \frac{r-1}{N-1}, \dots, \frac{r-(k-1)}{N-(k-1)}$$

are all greater than p_1 while the fractions

$$\frac{N-r}{N-k}, \frac{N-r-1}{N-k-1}, \dots, \frac{N-r-(m-1-k)}{N-(m-1)}$$

all exceed $\frac{N-r-(m-k)}{N-m}$ which equals $1 - p_2$. ■

When m is small relative to r and N , both fractions p_1 and p_2 are approximately equal to p . So this theorem justifies the earlier statement that sampling with and without replacement yield similar results when samples are small relative to the populations from which they were derived.

EXERCISES

Ex. 2.3.1. Suppose there are thirty balls in an urn, ten of which are black and the remaining twenty of which are red. Suppose three balls are selected from the urn (without replacement).

- (a) What is the probability that the sequence of draws is red-red-black?
- (b) What is the probability that the three draws result in exactly two red balls?

Ex. 2.3.2. This exercise explores how to use R to investigate the binomial approximation to the hypergeometric distribution.

- (a) A jar contains forty marbles – thirty white and ten black. Ten marbles are drawn at random from the jar. Use R to calculate the probability that exactly five of the marbles drawn are black. Do two separate computations, one under the assumption that the draws are with replacement and the other under the assumption that the draws are without replacement.
- (b) Repeat part (a) except now assume the jar contains 400 marbles – 300 white and 100 black.
- (c) Repeat part (a) except now assume the jar contains 4000 marbles – 3000 white and 1000 black.
- (d) Explain what you are observing with your results of parts (a), (b), and (c).

Ex. 2.3.3. Consider a room of one hundred people – forty men and sixty women.

- (a) If ten people are selected from the room, find the probability that exactly six are women. Calculate this probability with and without replacement and compare the decimal approximations of your two results.
- (b) If ten people are selected from the room, find the probability that exactly seven are women. Calculate this probability with and without replacement and compare the decimal approximations of your two results.
- (c) If 100 people are selected from the room, find the probability that exactly sixty are women. Calculate this probability with and without replacement and compare the two answers.
- (d) If 100 people are selected from the room, find the probability that exactly sixty-one are women. Calculate this probability with and without replacement and compare the two answers.

Ex. 2.3.4. Use the steps below to prove Theorem 2.3.2

- (a) Prove that $\frac{r!(N-k)!}{N!(r-k)!}$ equals

$$\frac{r}{N} \cdot \frac{r-1}{N-1} \cdots \frac{r-(k-1)}{N-(k-1)}.$$

- (b) Prove that $\frac{(N-r)!(N-m)!}{(N-k)!((N-r-(m-k))!)} \cdot \frac{N-r}{N-k} \cdot \frac{N-r-1}{N-k-1} \cdots \frac{N-r-(m-1-k)}{N-(m-1)}$ equals

$$\frac{N-r}{N-k} \cdot \frac{N-r-1}{N-k-1} \cdots \frac{N-r-(m-1-k)}{N-(m-1)}.$$

- (c) Use (a) and (b) to prove Theorem 2.3.2.

Ex. 2.3.5. A box contains W white balls and B black balls. A sample of n balls is drawn at random for some $n \leq \min(W, B)$. For $j = 1, 2, \dots, n$, let A_j denote the event that the ball drawn on the j^{th} draw is white. Let B_k denote the event that the sample of n balls contains exactly k white balls.

- (a) Find $P(A_j|B_k)$ if the sample is drawn with replacement.
- (b) Find $P(A_j|B_k)$ if the sample is drawn without replacement.

Ex. 2.3.6. For the problems below, assume a HyperGeo(N, r, m) distribution.

- (a) Calculate the ratio $\frac{P(\{k+1\})}{P(\{k\})}$.
(Assume that $\max\{0, m - (N - r)\} \leq k \leq \min\{r, m\}$ to avoid zero in the denominator).
- (b) Use (a) to calculate the mode of a HyperGeo(N, r, m).

Ex. 2.3.7. Biologists use a technique called “capture-recapture” to estimate the size of the population of a species that cannot be directly counted. The following exercise illustrates the role a hypergeometric distribution plays in such an estimate.

Suppose there is a species of unknown population size N . Suppose fifty members of the species are selected and given an identifying mark. Sometime later a sample of size twenty is taken from the population and it is found that four of the twenty were previously marked. The basic idea behind mark-recapture is that since the sample showed $\frac{4}{20} = 20\%$ marked members, that should also be a good estimate for the fraction of marked members of the species as a whole. However, for the whole species that fraction is $\frac{50}{N}$ which provides a population estimate of $N \approx 250$.

Looking more deeply at the problem, if the second sample is assumed to be done at random without replacement and with each member of the population equally likely to be selected, the resulting number of marked members should follow a HyperGeo($N, 50, 20$) distribution.

Under these assumptions use the formula for the mode calculated in the previous exercise to determine which values of N would cause a result of four marked members to be the most likely of the possible outcomes.

Ex. 2.3.8. The geometric distribution was first developed to determine the number of independent Bernoulli trials needed to observe the first success. When viewing the hypergeometric distribution as a series of dependent trials, the same question may be asked. Suppose we have a population of N people for which r have a certain characteristic and the remaining $N - r$ do not have that characteristic. Suppose an experiment consists of sampling (without replacement) repeatedly and recording the number of the sample that first corresponds to selecting someone with the specified characteristic. Answer the questions below.

- (a) What is S , the list of possible outcomes of this experiment?
- (b) For each $k \in S$, what is $P(\{k\})$?
- (c) Define $p = \frac{r}{N}$ and $p_1 = \frac{r-(k-1)}{N-(k-1)}$. Using the result from (b) prove the following bounds on the probability distribution:

$$p(1-p_1)^{k-1} \leq P(\{k\}) \leq p_1(1-p)^{k-1}$$

(As a consequence, when k is much smaller than r and N , the values of p_1 and p are approximately equal and the probabilities from (b) are closely approximated by a geometric distribution).

3

DISCRETE RANDOM VARIABLES

In the previous chapter many different distributions were developed out of Bernoulli trials. In that chapter we proceeded by creating new sample spaces for each new distribution, but when faced with many questions related to the same basic framework, it is usually clearer to maintain a single sample space and to define functions on that space whose outputs relate them to questions under consideration. Such functions are known as “random variables” and they will be the focus of this chapter.

3.1 RANDOM VARIABLES AS FUNCTIONS

EXAMPLE 3.1.1. Suppose a coin is flipped three times. Consider the probabilities associated with the following two questions:

- (a) How many coins will come up heads?
- (b) Which will be the first flip (if any) that shows heads?

At this point the answers to these questions should be easy to determine, but the purpose of this example is to emphasize how functions could be used to answer both within the context of a single sample space. Let S be a listing of all eight possible orderings of heads and tails on the three flips.

$$S = \{hhh, hht, hth, htt, thh, tht, tth, ttt\}$$

Now define two functions on S . Let X be the function that describes the total number of heads among the three flips and let Y be the function that describes the first flip that produces heads. So,

ω	$X(\omega)$	$Y(\omega)$
hhh	3	1
hht	2	1
hth	2	1
htt	1	1
thh	2	2
tht	1	2
tth	1	3
ttt	0	none

where $Y(ttt)$ is defined as “none” since there is no first time the coin produces heads.

Suppose we want to know the probability that exactly two of the three coins will be heads. The relevant event is $E = \{hht, hth, thh\}$, but in the pre-image notation of function theory this set may also be described as $X^{-1}(\{2\})$, the elements of S for which X produces an output of 2. This allows us to describe the probability of the event as:

$$P(\text{two heads}) = P(X^{-1}(\{2\})) = P(\{hht, hth, thh\}) = \frac{3}{8}$$

Rather than use the standard pre-image notation, it is more common in probability to write $(X = 2)$ for the set $X^{-1}(\{2\})$ since this emphasizes that we are considering outcomes for which the function X equals 2.

Similarly, if we wanted to know the probability that the first result of heads showed up on the third flip, that is a question that involves the function Y . Using the notation $(Y = 3)$ in place of $Y^{-1}(\{3\})$ the probability may be calculated as

$$P(\text{first heads on flip three}) = P(Y = 3) = P(\{\text{tth}\}) = \frac{1}{8}$$

As above we can compute the

$$P(X = 0) = \frac{1}{8}, P(X = 1) = \frac{3}{8}, \text{ and } P(X = 3) = \frac{1}{8}$$

and the

$$P(Y = 1) = \frac{1}{2}, P(Y = 2) = \frac{1}{4}, \text{ and } P(Y = \text{none}) = \frac{1}{8}.$$

Thus giving a complete description of how X and Y distribute the probabilities onto their range. In both cases only a single sample space was needed. Two different questions were approached by defining two different functions on that sample space. ■

The following theorem explains how the mechanism of the previous example may be more generally applied.

THEOREM 3.1.2. *Let S be a sample space with probability P and let $X : S \rightarrow T$ be a function. Then X generates a probability Q on T given by*

$$Q(B) = P(X^{-1}(B))$$

The probability Q is called the “distribution of X ” since it describes how X distributes the probability from S onto T .

The proof relies on two set-theoretic facts that we will take as given. The first is that $X^{-1}(\bigcup_{i=1}^{\infty} B_i) = \bigcup_{i=1}^{\infty} X^{-1}(B_i)$ and the second is the fact that if B_i and B_j are disjoint, then so are $X^{-1}(B_i)$ and $X^{-1}(B_j)$.

Proof - Let $B \subset T$. Since P is known to be a probability, $0 \leq P(X^{-1}(B)) \leq 1$, and so Q maps subsets of T into $[0, 1]$. Since X is a function into T , we know $X^{-1}(T) = S$. Therefore $Q(T) = P(X^{-1}(T)) = P(S) = 1$ and Q satisfies the first probability axiom.

To show Q satisfies the second axiom, suppose B_1, B_2, \dots are a countable collection of disjoint subsets of T .

$$\begin{aligned} Q\left(\bigcup_{i=1}^{\infty} B_i\right) &= P\left(X^{-1}\left(\bigcup_{i=1}^{\infty} B_i\right)\right) \\ &= P\left(\bigcup_{i=1}^{\infty} X^{-1}(B_i)\right) \\ &= \sum_{i=1}^{\infty} P(X^{-1}(B_i)) \\ &= \sum_{i=1}^{\infty} Q(B_i) \end{aligned}$$

As in the previous example, it is typical to write $(X \in B)$ in place of the notation $X^{-1}(B)$ to emphasize the fact that we are computing the probability that the function X takes a value in the set B . In practice, the new probability Q would rarely be used explicitly, but would be calculated in terms of the original probability P via the relationship described in the theorem.

EXAMPLE 3.1.3. A board game has a wheel that is to be spun periodically. The wheel can stop in one of ten equally likely spots. Four of these spots are red, three are blue, two are green, and one is black. Let X denote the color of the spot. Determine the distribution of X .

The function X is defined on a sample space S that consists of the ten spots the wheel could stop, and it takes values on the set of colors $T = \{\text{red}, \text{blue}, \text{green}, \text{black}\}$. Its distribution is a probability Q on the set of colors which can be determined by calculating the probability of each color.

For instance $Q(\{\text{red}\}) = P(X = \text{red}) = P(X^{-1}(\{\text{red}\})) = \frac{4}{10}$ since four of the ten spots on the wheel are red and all spots are equally likely. Similarly,

$$\begin{aligned} Q(\{\text{blue}\}) &= P(X = \text{blue}) = \frac{3}{10} \\ Q(\{\text{green}\}) &= P(X = \text{green}) = \frac{2}{10} \\ Q(\{\text{black}\}) &= P(X = \text{black}) = \frac{1}{10} \end{aligned}$$

completing the description of the distribution. ■

EXAMPLE 3.1.4. For a certain lottery, a three-digit number is randomly selected (from 000 to 999). If a ticket matches the number exactly, it is worth \$200. If the ticket matches exactly two of the three digits, it is worth \$20. Otherwise it is worth nothing. Let X be the value of the ticket. Find the distribution of X .

The function X is defined on $S = \{000, 001, \dots, 998, 999\}$ - the set of all one thousand possible three digit numbers. The function takes values on the set $\{0, 20, 200\}$, so the distribution Q is a probability on $T = \{0, 20, 200\}$.

First, $Q(\{200\}) = P(X = 200) = \frac{1}{1000}$ since only one of the one thousand three digit numbers is going to be an exact match.

Next, $Q(\{20\}) = P(X = 20)$, so it must be determined how many of the one thousand possibilities will have exactly two matches. There are $\binom{3}{2} = 3$ different ways to choose the two digits that will match. Those digits are determined at that point and the remaining digit must be one of the nine digits that do not match the third spot, so there are $3 \cdot 9 = 27$ three digit numbers that match exactly two digits. So $Q(\{20\}) = P(X = 20) = \frac{27}{1000}$.

Finally, since Q is a probability, $Q(\{0\}) = 1 - Q(\{20\}) - Q(\{200\}) = \frac{972}{1000}$. ■

It is frequently the case that we are interested in functions that have real-valued outputs and we reserve the term “random variable” for such a situation.

DEFINITION 3.1.5. A “discrete random variable” is a function $X : S \rightarrow T$ where S is a sample space equipped with a probability P , and T is a countable (or finite) subset of the real numbers.

From Theorem 3.1.2, P generates a probability on T and since it is a discrete space, the distribution may be determined by knowing the likelihood of each possible value of X . Because of this we define a function $f_X : T \rightarrow [0, 1]$ given by

$$f_X(t) = P(X = t)$$

referred to as a “probability mass function”. Then for any event $A \subset T$ the quantity $P(X \in A)$ may be computed via

$$P(X \in A) = \sum_{t \in A} f_X(t) = \sum_{t \in A} P(X = t).$$

The function from Example 3.1.4 is a discrete random variable because it takes on one of the real values 0, 20, or 200. We calculated its probability mass function when describing its distribution and it is given by

$$f_X(0) = \frac{972}{1000}, f_X(20) = \frac{27}{1000}, f_X(200) = \frac{1}{1000}.$$

The function from Example 3.1.3 is not a discrete random variable by the above definition because its range is a collection of colors, not real numbers.

3.1.1 Common Distributions

When studying random variables it is often more important to know how they distribute probability onto their range than how they actually act as functions on their domains. As such it is useful to have a notation that recognizes the fact that two functions may be very different in terms of where they map domain elements, but nevertheless have the same range and produce the same distribution on this range.

DEFINITION 3.1.6. Let $X : S \rightarrow T$ and $Y : S \rightarrow T$ be discrete random variables. We say X and Y have equal distribution provided $P(X = t) = P(Y = t)$ for all $t \in T$.

There are many distributions which appear frequently enough they deserve their own special names for easy identification. We shall use the symbol \sim to mean “is distributed as” or “is equal in distribution to”. For example, in the definition below $X \sim \text{Bernoulli}(p)$ should be read as “ X has a Bernoulli(p) distribution”. This says nothing explicit about how X behaves as a function on its domain, but completely describes how X distributes probability onto its range.

DEFINITION 3.1.7. The following are common discrete distributions which we have seen arise previously in the text.

- (a) $X \sim \text{Uniform}(\{1, 2, \dots, n\})$: Let $n \geq 1$ be an integer. If X is a random variable such that $P(X = k) = \frac{1}{n}$ for all $1 \leq k \leq n$ then we say that X is a uniform random variable on the set $\{1, 2, \dots, n\}$.
- (b) $X \sim \text{Bernoulli}(p)$: Let $0 \leq p \leq 1$. When X is a random variable such that $P(X = 1) = p$ and $P(X = 0) = 1 - p$ we say that X is a Bernoulli random variable with parameter p . This takes the concept of a “Bernoulli trial” which we have previously discussed and puts it in the context of a random variable where 1 corresponds to success and 0 corresponds to failure.

- (c) $X \sim \text{Binomial}(n, p)$: Let $0 \leq p \leq 1$ and let $n \geq 1$ be an integer. If X is a random variable taking values in $\{0, 1, \dots, n\}$ having a probability mass function

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

for all $0 \leq k \leq n$, then X is a binomial random variable with parameters n and p . We have seen that such a quantity describes the number of successes in n Bernoulli trials.

- (d) $X \sim \text{Geometric}(p)$: Let $0 < p < 1$. If X is a random variable with values in $\{1, 2, 3, \dots\}$ and a probability mass function

$$P(X = k) = p \cdot (1-p)^{k-1}$$

for all $k \geq 1$, then X is a geometric random variable with parameter p . Such a random variable arises when determining how many Bernoulli trials must be attempted before seeing the first success.

- (e) $X \sim \text{Negative Binomial } (r, p)$: Let $0 < p < 1$. If X is a random variable with values in $\{r, r+1, r+2, \dots\}$ and a probability mass function

$$P(X = k) = \binom{k-1}{r-1} p^r \cdot (1-p)^{k-r}$$

for all $k \geq r$, then X is a negative binomial random variable with parameters (r, p) . Such a random variable arises when determining how many Bernoulli trials must be attempted before seeing r successes.

- (f) $X \sim \text{Poisson}(\lambda)$: Let $\lambda > 0$. When X is a random variable with values in $\{0, 1, 2, \dots\}$ such that its probability mass function is

$$P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

for all $k \geq 0$, then X is called a Poisson random variable with parameter λ . We first used these distributions as approximations to a $\text{Binomial}(n, p)$ when n was large and p was small.

- (g) $X \sim \text{HyperGeo}(N, r, m)$: Let N , r , and m be positive integers for which $r < N$ and $m < N$. Let X be a random variable taking values in the integers between $\min\{m, r\}$ and $\max\{0, m - (N - r)\}$ inclusive with probability mass function

$$P(X = k) = \frac{\binom{r}{k} \binom{N-r}{m-k}}{\binom{N}{m}}$$

The random variable X is called hypergeometric with parameters N , r , and m . Such quantities occur when sampling without replacement.

EXERCISES

Ex. 3.1.1. Consider the experiment of flipping a coin four times and recording the sequence of heads and tails. Let S be the sample space of all sixteen possible orderings of the results. Let X be the function on S describing the number of tails among the flips. Let Y be the function on S describing the first flip (if any) to come up tails.

- (a) Create a table as in Example 3.1.1 describing functions X and Y .
- (b) Use the table to calculate $P(X = 2)$.
- (c) Use the table to calculate $P(Y = 3)$.

Ex. 3.1.2. A pair of fair dice are thrown. Let X represent the larger of the two values on the dice and let Y represent the smaller of the two values.

- (a) Describe S , the domain of functions X and Y . How many elements are in S ?
- (b) What are the ranges of X and Y . Do X and Y have the same range? Why or why not?
- (c) Describe the distribution of X and describe the distribution of Y by finding the probability mass function of each. Is it true that X and Y have the same distribution ?

Ex. 3.1.3. A pair of fair dice are thrown. Let X represent the number of the first die and let Y represent the number of the second die.

- (a) Describe S , the domain of functions X and Y . How many elements are in S ?
- (b) Describe T , the range of functions X and Y . How many elements are in T ?
- (c) Describe the distribution of X and describe the distribution of Y by finding the probability mass function of each. Is it true that X and Y have the same distribution ?
- (d) Are X and Y the same function? Why or why not?

Ex. 3.1.4. Use the \sim notation to classify the distributions of the random variables described by the scenarios below. For instance, if a scenario said, “let X be the number of heads in three flips of a coin” the appropriate answer would be $X \sim \text{Binomial}(3, \frac{1}{2})$ since that describes the number of successes in three Bernoulli trials.

- (a) Let X be the number of 5's seen in four die rolls. What is the distribution of X ?
- (b) Each ticket in a certain lottery has a 20% chance to be a prize-winning ticket. Let Y be the number of tickets that need to be purchased before seeing the first prize-winner. What is the distribution of Y ?
- (c) A class of ten students is comprised of seven women and three men. Four students are randomly selected from the class. Let Z denote the number of men among the four randomly selected students. What is the distribution of Z ?

Ex. 3.1.5. Suppose X and Y are random variables.

- (a) Explain why $X + Y$ is a random variable.
- (b) Theorem 3.1.2 does not require that X be real-valued. Why do you suppose that our definition of “random variable” insisted that such functions should be real-valued?

Ex. 3.1.6. Let $X : S \rightarrow T$ be a discrete random variable. Suppose $\{B_i\}_{i \geq 1}$ are sequence of events in T then show that $X^{-1}(\bigcup_{i=1}^{\infty} B_i) = \bigcup_{i=1}^{\infty} X^{-1}(B_i)$ and that if B_i and B_j are disjoint, then so are $X^{-1}(B_i)$ and $X^{-1}(B_j)$.

3.2 INDEPENDENT AND DEPENDENT VARIABLES

Most interesting problems require the consideration of several different random variables and an analysis of the relationships among them. We have already discussed what it means for a collection of events to be independent and it is useful to extend this notion to random variables as well. As with events we will first describe the notion of pairwise independence of two objects before defining mutual independence of an arbitrary connection of objects.

3.2.1 Independent Variables

DEFINITION 3.2.1. (Independence of a Pair of Random Variables) Two random variables X and Y are independent if $(X \in A)$ and $(Y \in B)$ are independent for every event A in the range of X and every event B in the range of Y .

As events become more complicated and involve multiple random variables, a notational shorthand will become useful. It is common in probability to write $(X \in A, Y \in B)$ for the event $(X \in A) \cap (Y \in B)$ and we will begin using this convention at this point.

Further, even though the definition of $X : S \rightarrow T$ and $Y : S \rightarrow U$ being independent random variables requires that $(X \in A)$ and $(Y \in B)$ be independent for all events $A \subset T$ and $B \subset U$, for discrete random variables it is enough to verify the events $(X = t)$ and $(Y = u)$ are independent events for all $t \in T$ and $u \in U$ to conclude they are independent (See Exercise 3.2.12).

EXAMPLE 3.2.2. When we originally considered the example of rolling a pair of dice, we viewed the results as thirty-six equally likely outcomes. However, it is also possible to view the result of each die as a random variable in its own right, and then consider the possible results of the pair of random variables. Let $X, Y \sim \text{Uniform}(\{1, 2, 3, 4, 5, 6\})$ and suppose X and Y are independent. If $x, y \in \{1, 2, 3, 4, 5, 6\}$ what is $P(X = x, Y = y)$?

By independence $P(X = x, Y = y) = P(X = x)P(Y = y) = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}$. Therefore, the result is identical to the original perspective – each of the thirty-six outcomes of the pair of dice is equally likely. ■

DEFINITION 3.2.3. (Mutual Independence of Random Variables) A finite collection of random variables X_1, X_2, \dots, X_n is mutually independent if the sets $(X_j \in A_j)$ are mutually independent for all events A_j in the ranges of the corresponding X_j .

An arbitrary collection of random variables X_t where $t \in I$ for some index set I is mutually independent if every finite sub-collection is mutually independent.

For many problems it is useful to think about repeating a single experiment many times with the results of each repetition being independent from every other. Though the results are assumed

to be independent, the experiment itself remains the same, so the random variables produced all have the same distribution. The resulting sequence of random variables X_1, X_2, X_3, \dots is referred to as “i.i.d.” (standing for “independent and identically distributed”). When considering such sequences we will sometimes write X_1, X_2, X_3, \dots are i.i.d. with distribution X , where X is a random variable that shares their common distribution.

EXAMPLE 3.2.4. Let X_1, \dots, X_n be i.i.d. with a $\text{Geometric}(p)$ distribution. What is the probability that all of these random variables are larger than some positive integer j ?

As a preliminary calculation, if $X \sim \text{Geometric}(p)$ and if $j \geq 1$ is an integer we may determine $P(X > j)$.

$$\begin{aligned} P(X > j) &= \sum_{i=j+1}^{\infty} P(X = i) \\ &= \sum_{i=j+1}^{\infty} p(1-p)^{i-1} \\ &= \frac{p \cdot (1-p)^j}{1 - (1-p)} \\ &= (1-p)^j \end{aligned}$$

But each of X_1, X_2, \dots, X_n have this distribution, so using the computation above, together with independence,

$$\begin{aligned} P(X_1 > j, X_2 > j, \dots, X_n > j) &= P(X_1 > j)P(X_2 > j) \dots P(X_n > j) \\ &= (1-p)^j \cdot (1-p)^j \dots \cdot (1-p)^j \\ &= (1-p)^{nj} \end{aligned}$$



3.2.2 Conditional, Joint, and Marginal Distributions

Consider a problem involving two random variables. Let X be the number of centimeters of rainfall in a certain forest in a given year, and let Y be the number of square meters of the forest burned by fires that same year. It seems these variables should be related; knowing one should affect the probabilities associated with the values of the other. Such random variables are not independent of each other and we now introduce several ways to compute probabilities under such circumstances. An important concept toward this end is the notion of a “conditional distribution” which reflects the fact that the occurrence of an event may affect the likely values of a random variable.

DEFINITION 3.2.5. Let X be a random variable on a sample space S and let $A \subset S$ be an event such that $P(A) > 0$. Then the probability Q described by

$$Q(B) = P(X \in B | A) \tag{3.2.1}$$

is called the “conditional distribution” of X given the event A .

As with any discrete random variable, the distribution is completely determined by the probabilities associated with each possible value the random variable may assume. This means the conditional distribution may be considered known provided the values of $P(X = a|A)$ are known for every $a \in \text{Range}(X)$. Though this definition allows for A to be any sort of event, in this section we will mainly consider examples where A describes the outcome of some random variable. So a notation like $P(X|Y = b)$ will be the conditional distribution of the random variable X given that the variable Y is known to have the value b .

In many cases random variables are dependent in such a way that the distribution of one variable is known in terms of the values taken on by another.

EXAMPLE 3.2.6. Let $X \sim \text{Uniform}(\{1, 2\})$ and let Y be the number of heads in X tosses of a fair coin. Clearly X and Y should not be independent. In particular, a result of $Y = 0$ could occur regardless of the value of X , but a result of $Y = 2$ guarantees that $X = 2$ since two heads could not be observed with just one flip on the coin. Any information regarding X or Y may influence the distribution of the other, but the description of the variables makes it clearest how Y depends on X . If $X = 1$ then Y is the number of heads in one flip of a fair coin. Letting A be the event ($X = 1$) and using the terminology of (3.2.1) from Definition 3.2.5, we can say the conditional distribution of Y given that $X = 1$ is a Bernoulli($\frac{1}{2}$). We will use the notation

$$(Y|X = 1) \sim \text{Bernoulli}\left(\frac{1}{2}\right)$$

to indicate this fact. In other words, this notation means the same thing as the pair of equations

$$\begin{aligned} P(Y = 0|X = 1) &= \frac{1}{2} \\ P(Y = 1|X = 1) &= \frac{1}{2} \end{aligned}$$

If $X = 2$ then Y is the number of heads in two flips of a fair coin and therefore $(Y|X = 2) \sim \text{Binomial}(2, \frac{1}{2})$ which means the following three equations hold:

$$\begin{aligned} P(Y = 0|X = 2) &= \frac{1}{4} \\ P(Y = 1|X = 2) &= \frac{1}{2} \\ P(Y = 2|X = 2) &= \frac{1}{4} \end{aligned}$$

■

The conditional probabilities of the previous example were easily determined in part because the description of Y was already given in terms of X , but frequently random variables may be dependent in some way that is not so explicitly described. A more general method of expressing the dependence of two (or more) variables is to present the probabilities associated with all combinations of possible values for every variable. This is known as their joint distribution.

DEFINITION 3.2.7. If X and Y are discrete random variables, the “joint distribution” of X and Y is the probability Q on pairs of values in the ranges of X and Y defined by

$$Q((a, b)) = P(X = a, Y = b).$$

The definition may be expanded to a finite collection of discrete random variables X_1, X_2, \dots, X_n for which the joint distribution of all n variables is the probability defined by

$$Q((a_1, a_2, \dots, a_n)) = P(X_1 = a_1, X_2 = a_2, \dots, X_n = a_n).$$

In the above definition as discussed before for any event D ,

$$Q(D) = \sum_{(a_1, a_2, \dots, a_n) \in D} Q((a_1, a_2, \dots, a_n)).$$

For a pair of random variables with few possible outcomes, it is common to describe the joint distribution using a chart for which the columns correspond to possible X values, the rows to possible Y values, and for which the entries of the chart are probabilities.

EXAMPLE 3.2.8. Let X and Y be the dependent variables described in Example 3.2.6. The X variable will be either 1 or 2. The Y variable could be as low as 0 (if no heads are flipped) or as high as 2 (if two coins are flipped and both show heads). Since $\text{Range}(X) = \{1, 2\}$ and since $\text{Range}(Y) = \{0, 1, 2\}$, the pair (X, Y) could potentially be any of the six possible pairings (though, in fact, one of the pairings has probability zero). To find the joint distribution of X and Y we must calculate the probabilities of each possibility. In this case the values may be obtained using the definition of conditional probability. For instance,

$$P(X = 1, Y = 0) = P(Y = 0|X = 1) \cdot P(X = 1) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

and

$$P(X = 1, Y = 2) = P(Y = 2|X = 1) \cdot P(X = 1) = 0 \cdot \frac{1}{2} = 0$$

The entire joint distribution $P(X = a, Y = b)$ is described by the following chart.

	$X = 1$	$X = 2$
$Y = 0$	1/4	1/8
$Y = 1$	1/4	1/4
$Y = 2$	0	1/8



Knowing the joint distribution of random variables gives a complete picture of the probabilities associated with those variables. From that information it is possible to compute all conditional probabilities of one variable from another. For instance, in the example analyzed above, the variable Y was originally described in terms of how it depended on X . However, this also means that X should be dependent on Y . The joint distribution may be used to determine how.

EXAMPLE 3.2.9. Let X and Y be the variables of Example 3.2.6. How may the conditional distributions of X given values of Y be determined?

There will be three different conditional distributions depending on whether $Y = 0$, $Y = 1$, or $Y = 2$. Below we will solve the $Y = 0$ case. The other two cases will be left as exercises. The

conditional distribution of X given that $Y = 0$ is determined by the values of $P(X = 1|Y = 0)$ and $P(X = 2|Y = 0)$ both of which may be computed using Bayes' rule.

$$\begin{aligned} P(X = 1|Y = 0) &= \frac{P(Y = 0|X = 1) \cdot P(X = 1)}{P(Y = 0)} \\ &= \frac{P(Y = 0|X = 1) \cdot P(X = 1)}{P(Y = 0|X = 1) \cdot P(X = 1) + P(Y = 0|X = 2) \cdot P(X = 2)} \\ &= \frac{(1/2)(1/2)}{(1/2)(1/2) + (1/4)(1/2)} = \frac{2}{3} \end{aligned}$$

Since the only values for X are 1 and 2 it must be that $P(X = 2|Y = 0) = \frac{1}{3}$. ■

Just because X and Y are dependent on each other doesn't mean they need to be thought of as a pair. It still makes sense to talk about the distribution of X as a random variable in its own right while ignoring its dependence on the variable Y . When there are two or more variables under discussion, the distribution of X alone is sometimes called the "marginal" distribution of X because it can be computed using the margins of the chart describing the joint distribution of X and Y .

EXAMPLE 3.2.10. Continue with X and Y as described in Example 3.2.6. Below is the chart describing the joint distribution of X and Y that was created in Example 3.2.8, but with the addition of one column on the right and one row at the bottom. The entries in the extra column are the sums of the values in the corresponding row; likewise the entries in the extra row are the sums of the values in the corresponding column.

	$X = 1$	$X = 2$	Sum
$Y = 0$	1/4	1/8	3/8
$Y = 1$	1/4	1/4	4/8
$Y = 2$	0	1/8	1/8
Sum	1/2	1/2	

The values in the right hand margin (column) exactly describe the distribution of Y . For instance the event $(Y = 0)$ can be partitioned into two disjoint events $(X = 1, Y = 0) \cup (X = 2, Y = 0)$ each of which is already described in the joint distribution chart. Adding them together gives the result that $P(Y = 0) = \frac{3}{8}$. In a similar fashion, the bottom margin (row) describes the distribution of X . This extended chart also makes it numerically clearer why these two random variables are dependent. For instance,

$$P(X = 1, Y = 0) = \frac{1}{4} \quad \text{while} \quad P(X = 1) \cdot P(Y = 0) = \frac{3}{16}$$

Since these quantities are unequal, the random variables cannot be independent. ■

In general, knowing the marginal distributions of X and Y is not sufficient information to reconstruct their joint distribution. This is because the marginal distributions do not provide any information about how the random variables relate to each other. However, if X and Y happen to be independent, then their joint distribution may easily be computed from the marginals since

$$P(X = x, Y = y) = P(X = x)P(Y = y)$$

3.2.3 Memoryless Property of the Geometric Random Variable

It is also possible to calculate conditional probabilities of a random variable based on subsets of its own values. A particularly important example of this is the "memoryless property" of geometric random variables.

EXAMPLE 3.2.11. Suppose we toss a fair coin until the first head appears. Let X be the number of tosses performed. We have seen in Example 2.1.2 that $X \sim \text{Geometric}(\frac{1}{2})$. Note that if m is a positive integer,

$$P(X > m) = \sum_{k=m+1}^{\infty} P(X = k) = \sum_{k=m+1}^{\infty} \frac{1}{2^k} = \frac{1}{2^m}$$

Now let n be a positive integer and suppose we take the event $(X > n)$ as given. In other words, we assume we know that none of the first n flips resulted in heads. What is the conditional distribution of X given this new information? A routine calculation shows

$$P(X > n+m | X > n) = \frac{P(X > n+m)}{P(X > n)} = \frac{\frac{1}{2^{n+m}}}{\frac{1}{2^n}} = \frac{1}{2^m}$$

As a consequence,

$$P(X > n+m | X > n) = P(X > m). \quad (3.2.2)$$

Given that a result of heads has not occurred by the n^{th} flip, the probability that such a result will require at least m more flips is identical to the (non-conditional) probability the result would have required more than m flips from the start. In other words, if we know that the first n flips have not yet produced a head, the number of additional flips required to observe the first head still is a Geometric($\frac{1}{2}$) random variable. This is called the “memoryless property” of the geometric distribution since it can be interpreted to mean that when waiting times are geometrically distributed, no matter how long we wait for an event to occur, the future waiting time always looks the same given that the event has not occurred yet. The result remains true of geometric variables of any parameter p , a fact which we leave as an exercise. ■

3.2.4 Multinomial Distributions

Consider a situation similar to that of Bernoulli trials, but instead of results of each attempt limited to success or failure, suppose there are many different possible results for each trial. As with the Bernoulli trial cases we assume that the trials are mutually independent, but identically distributed. In the next example we will show how to calculate the joint distribution for the random variables representing the number of times each outcome occurs.

EXAMPLE 3.2.12. Suppose we perform n i.i.d. trials each of which has k different possible outcomes. For $j = 1, 2, \dots, k$, let p_j represent the probability any given trial results in the j^{th} outcome and let X_j represent the number of the n trials that result in the j^{th} outcome. The joint distribution of all of the random variables X_1, X_2, \dots, X_k is called a “multinomial distribution”.

Let $B(x_1, x_2, \dots, x_k) = \{X_1 = x_1, X_2 = x_2, \dots, X_k = x_k\}$. Then,

$$P(B(x_1, x_2, \dots, x_k)) = \sum_{\omega \in B(x_1, x_2, \dots, x_k)} P(\{\omega\})$$

Each $\omega \in B(x_1, x_2, \dots, x_k)$ is an element in the sample space corresponding to the j^{th} outcome occurring exactly x_j times. Since the trials are independent, and since an outcome j occurs in x_j trials, each of which had probability p_j , this means

$$P(\{\omega\}) = \prod_{j=1}^k p_j^{x_j}$$

Consequently, each outcome in $B(x_1, x_2, \dots, x_k)$ has the same probability. So to determine the likelihood of the event, we need only determine $|B(x_1, x_2, \dots, x_k)|$, the number of outcomes the event contains. The calculation of this quantity is a combinatorial problem; it is the number of ways of allocating n balls in k boxes, such that x_j of them fall into box j . We leave it as an exercise to prove that

$$|B(x_1, x_2, \dots, x_k)| = \frac{n!}{x_1! x_2! \dots x_k!}$$

With that computation complete, the joint distribution of X_1, X_2, \dots, X_k is given by

$$P(X_1 = x_1, \dots, X_k = x_k) = \begin{cases} \frac{n!}{x_1! x_2! \dots x_k!} \prod_{j=1}^k p_j^{x_j} & \text{if } x_j \in \{0, 1, \dots, n\} \\ & \text{and } \sum_{j=1}^k x_j = n \\ 0 & \text{otherwise} \end{cases}$$

■

EXERCISES

Ex. 3.2.1. An urn has four balls labeled 1, 2, 3, and 4. A first ball is drawn and its number is denoted by X . A second ball is then drawn from the three remaining balls in the urn and its number is denoted by Y .

- (a) Calculate $P(X = 1)$.
- (b) Calculate $P(Y = 2|X = 1)$.
- (c) Calculate $P(Y = 2)$.
- (d) Calculate $P(X = 1, Y = 2)$.
- (e) Are X and Y independent? Why or why not?

Ex. 3.2.2. Two dice are rolled. Let X denote the sum of the dice and let Y denote the value of the first die.

- (a) Calculate $P(X = 7)$ and $P(Y = 4)$.
- (b) Calculate $P(X = 7, Y = 4)$.
- (c) Calculate $P(X = 5)$ and $P(Y = 4)$.
- (d) Calculate $P(X = 5, Y = 4)$.
- (e) Are X and Y independent? Why or why not?

Ex. 3.2.3. Let X and Y be the variables described in Example 3.2.6.

- (a) Determine the conditional distribution of X given that $Y = 1$.
- (b) Determine the conditional distribution of X given that $Y = 2$.

Ex. 3.2.4. Let X and Y be random variables with joint distribution given by the chart below.

	$X = 0$	$X = 1$	$X = 2$
$Y = 0$	1/12	0	3/12
$Y = 1$	2/12	1/12	0
$Y = 2$	3/12	1/12	1/12

- (a) Compute the marginal distributions of X and Y .
- (b) Compute the conditional distribution of X given that $Y = 2$.
- (c) Compute the conditional distribution of Y given that $X = 2$.
- (d) Carry out a computation to show that X and Y are not independent.

Ex. 3.2.5. Let X be a random variable with range $\{0, 1\}$ and distribution

$$P(X = 0) = \frac{1}{3} \quad \text{and} \quad P(X = 1) = \frac{2}{3}$$

and let Y be a random variable with range $\{0, 1, 2\}$ and distribution

$$P(Y = 0) = \frac{1}{5}, \quad P(Y = 1) = \frac{1}{5}, \quad \text{and} \quad P(Y = 2) = \frac{3}{5}$$

Suppose that X and Y are independent. Create a chart describing the joint distribution of X and Y .

Ex. 3.2.6. Consider six independent trials each of which are equally likely to produce a result of 1, 2, or 3. Let X_j denote the number of trials that result in j . Calculate $P(X_1 = 1, X_2 = 2, X_3 = 3)$.

Ex. 3.2.7. Prove the combinatorial fact from Example 3.2.12 in the following way. Let $A_n(x_1, x_2, \dots, x_k)$ denote the number of ways of putting n balls into k boxes in such a way that exactly x_j balls wind up in box j for $j = 1, 2, \dots, k$.

- (a) Prove that $A_n(x_1, x_2, \dots, x_k) = \binom{n}{x_1} A_{n-x_1}(x_2, x_3, \dots, x_k)$.
- (b) Use part (a) and induction to prove that $A_n(x_1, x_2, \dots, x_k) = \frac{n!}{x_1! x_2! \dots x_k!}$.

Ex. 3.2.8. Let X be the result of a fair die roll and let Y be the number of heads in X coin flips.

- (a) Both X and $(Y|X = n)$ can be written in terms of common distributions using the \sim notation. What is the distribution of X ? What is the distribution of $(Y|X = n)$ for $n = 1, \dots, 6$?
- (b) Determine the joint distribution for X and Y .
- (c) Calculate the marginal distribution of Y .
- (d) Compute the conditional distribution of X given that $Y = 6$.
- (e) Compute the conditional distribution of X given that $Y = 0$.
- (f) Perform a computation to prove that X and Y are not independent.

Ex. 3.2.9. Suppose the number of earthquakes that occur in a year, anywhere in the world, is a Poisson random variable with mean λ . Suppose the probability that any given earthquake has magnitude at least 5 on the Richter scale is p independent of all other quakes. Let $N \sim \text{Poisson}(\lambda)$ be the number of earthquakes in a year and let M be the number of earthquakes in a year with magnitude at least 5, so that $(M|N = n) \sim \text{Binomial}(n, p)$.

(a) Calculate the joint distribution of M and N .

(b) Show that the marginal distribution of M is determined by

$$P(M = m) = \frac{1}{m!} e^{-\lambda} (\lambda p)^m \sum_{n=m}^{\infty} \frac{\lambda^{n-m}}{(n-m)!} (1-p)^{n-m}$$

for $m > 0$.

(c) Perform a change of variables (where $k = n - m$) in the infinite series from part (b) to prove

$$P(M = m) = \frac{1}{m!} e^{-\lambda} (\lambda p)^m \sum_{k=0}^{\infty} \frac{(\lambda(1-p))^k}{k!}$$

(d) Use part (c) together with the infinite series equality $e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$ to conclude that $M \sim \text{Poisson}(\lambda p)$.

Ex. 3.2.10. Let X be a discrete random variable which has $\mathbb{N} = \{1, 2, 3, \dots\}$ as its range. Suppose that for all positive integers m and n , X has the memoryless property – $P(X > n + m | X > n) = P(X > m)$. Prove that X must be a geometric random variable. [Hint: Define $p = P(X = 1)$ and use the memoryless property to calculate $P(X = n)$ inductively].

Ex. 3.2.11. A discrete random variable X is called “constant” if there is a single value c for which $P(X = c) = 1$.

- (a) Prove that if X is a constant discrete random variable then X is independent of itself.
- (b) Prove that if X is a discrete random variable which is independent of itself, then X must be constant. [Hint: It may help to look at Exercise 1.4.8].

Ex. 3.2.12. Let $X : S \rightarrow T$ and $Y : S \rightarrow U$ be discrete random variables. Show that if

$$P(X = t, Y = u) = P(X = t)P(Y = u)$$

for all $t \in T$ and $u \in U$ then X and Y are independent random variables.

3.3 FUNCTIONS OF RANDOM VARIABLES

There are many circumstances where we want to consider functions applied to random variables as inputs of functions. For a simple geometric example, suppose a rectangle is selected in such a way that its width X and its length Y are both random variables with known joint distribution. The area of the rectangle is $A = XY$ and since X and Y are random, it should be that A is random as well. How may the distribution of A be calculated from the joint distribution of X and Y ? In general, if a new random variable Z depends on random variables X_1, X_2, \dots, X_n which have a given joint distribution, how may the distribution of Z be calculated from what is already known? In this section we discuss the answers to such questions and also address related issues surrounding independence.

If $X : S \rightarrow T$ is a random variable and if $f : T \rightarrow R$ is a function, then the quantity $f(X)$ makes sense as a composition of functions $f \circ X : S \rightarrow R$. In fact, since $f(X)$ is defined on the sample space S , this new composition is itself a random variable.

The same reasoning holds for functions of more than one variable. If X_1, X_2, \dots, X_n are random variables then $f(X_1, X_2, \dots, X_n)$ is a random variable provided f is defined for the values the X_j variables produce. Below we illustrate how to calculate the distribution of $f(X_1, X_2, \dots, X_n)$ in terms of the joint distribution of the X_j input variables. We demonstrate the method with several examples followed by a general theorem.

3.3.1 Distribution of $f(X)$ and $f(X_1, X_2, \dots, X_n)$

The distribution of $f(X)$ involves the probability of events such as $(f(X) = a)$ for values of a that the function may produce. The key to calculating this probability is that these events may be rewritten in terms of the input values of X instead of the output values of $f(X)$.

EXAMPLE 3.3.1. Let $X \sim \text{Uniform}(\{-2, -1, 0, 1, 2\})$ and let $f(x) = x^2$. Determine the range and distribution of $f(X)$.

Since $f(X) = X^2$, the values that $f(X)$ produces are the squares of the values that X produces. Squaring the values in $\{-2, -1, 0, 1, 2\}$ shows the range of $f(X)$ is $\{0, 1, 4\}$. The probabilities that $f(X)$ takes on each of these three values determine the distribution of $f(X)$ and these probabilities can be easily calculated from the known probabilities associated with X .

$$\begin{aligned} P(f(X) = 0) &= P(X = 0) = \frac{1}{5} \\ P(f(X) = 1) &= P((X = 1) \cup (X = -1)) = \frac{1}{5} + \frac{1}{5} = \frac{2}{5} \\ P(f(X) = 4) &= P((X = 2) \cup (X = -2)) = \frac{1}{5} + \frac{1}{5} = \frac{2}{5} \end{aligned}$$

■

A complication with this method is that there may be many different inputs that produce the same output. Sometimes a problem requires careful consideration of all ways that a given output may be produced. For instance,

EXAMPLE 3.3.2. What is the probability the sum of three dice will equal six? Let X , Y , and Z be the results of the first, second, and third die respectively. These are i.i.d. random variables each distributed as $\text{Uniform}(\{1, 2, 3, 4, 5, 6\})$. A sum of six can be arrived at in three distinct ways:

- Case I: through three rolls of 2;
- Case II: through one roll of 3, one roll of 2, and one roll of 1; or
- Case III: through one roll of 4 and two rolls of 1

The first of these is the simplest to deal with since independence gives

$$P(X = 2, Y = 2, Z = 2) = P(X = 2) \cdot P(Y = 2) \cdot P(Z = 2) = \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{216}$$

The other cases involve a similar computation, but are complicated by the consideration of which number shows up on which die. For instance, both events $(X = 1, Y = 2, Z = 3)$ and $(X = 3, Y = 2, Z = 1)$ are included as part of Case II as are four other permutations of the numbers. Likewise Case III includes three permutations, one of which is $(X = 4, Y = 1, Z = 1)$. Putting all three cases together,

$$\begin{aligned}
 P(\text{sum of } 6) &= P(\text{Case I}) + P(\text{Case II}) + P(\text{Case III}) \\
 &= \frac{1}{216} + 6\left(\frac{1}{216}\right) + 3\left(\frac{1}{216}\right) \\
 &= \frac{5}{108}
 \end{aligned}$$

So there is slightly less than a 5% chance three rolled dice will produce a sum of six. ■

This method may also be used to show relationships among the common (named) distributions that have been previously described, as in the next two examples.

EXAMPLE 3.3.3. Let $X, Y \sim \text{Bernoulli}(p)$ be two independent random variables. If $Z = X + Y$, show that $Z \sim \text{Binomial}(2, p)$.

This result should not be surprising given how Bernoulli and binomial distributions arose in the first place. Each of X and Y produces a value of 0 if the corresponding Bernoulli trial was a failure and 1 if the trial was a success. Therefore $Z = X + Y$ equals the total number of successes in two independent Bernoulli trials, which is exactly what led us to the binomial distribution in the first place. However, it is instructive to consider how this problem relates to the current topic of discussion.

Since each of X and Y is either 0 or 1 the possible values of Z are in the set $\{0, 1, 2\}$. A result of $Z = 0$ can only occur if both X and Y are zero. So,

$$\begin{aligned}
 P(Z = 0) &= P(X = 0, Y = 0) \\
 &= P(X = 0) \cdot P(Y = 0) \\
 &= (1 - p)(1 - p) = (1 - p)^2
 \end{aligned}$$

Similarly, $P(Z = 2) = P(X = 1) \cdot P(Y = 1) = p^2$.

There are two different ways that Z could equal 1, either $X = 1$ and $Y = 0$, or $X = 0$ and $Y = 1$. So,

$$\begin{aligned}
 P(Z = 1) &= P((X = 1, Y = 0) \cup (X = 0, Y = 1)) \\
 &= P(X = 1, Y = 0) + P(X = 0, Y = 1) \\
 &= p(1 - p) + (1 - p)p = 2p(1 - p)
 \end{aligned}$$

These values of $P(Z = 0)$, $P(Z = 1)$, and $P(Z = 2)$ are exactly what define $Z \sim \text{Binomial}(2, p)$. ■

Two of the previous three examples involving adding random variables together. In fact, addition is one of the most common examples of applying functions to random quantities. In the previous situations, calculating the distribution of the sum was relatively simple because the component variables only had finitely many outcomes. But now suppose X and Y are random variables taking values in $\{0, 1, 2, \dots\}$ and suppose $Z = X + Y$. How could $P(Z = n)$ be calculated?

Since both X and Y are non-negative and since $Z = X + Y$, the value of Z must be at least as large as either X or Y individually. If $Z = n$, then X could take on any value $j \in \{0, 1, \dots, n\}$, but once that value is determined, the value of Y is compelled to be $n - j$ to give the appropriate sum. In other words, the event $(Z = n)$ partitions into the following union.

$$(Z = n) = \bigcup_{j=0}^n (X = j, Y = n - j)$$

When X and Y are independent, this means

$$\begin{aligned} P(Z = n) &= P\left(\bigcup_{j=0}^n (X = j, Y = n - j)\right) \\ &= \sum_{j=0}^n P(X = j, Y = n - j) \\ &= \sum_{j=0}^n P(X = j) \cdot P(Y = n - j) \end{aligned}$$

Such a computation is usually referred to as a “convolution” which will be addressed more generally later in the text. It occurs regularly when determining the distribution of sums of independent random variables.

EXAMPLE 3.3.4. Let $X \sim \text{Poisson}(\lambda_1)$ and $Y \sim \text{Poisson}(\lambda_2)$ be independent random variables.

- (a) Let $Z = X + Y$. Find the distribution of Z .
- (b) Find the conditional distribution of $X | Z$.

Now

$$P(X = x, Y = y) = P(X = x) \cdot P(Y = y) = e^{-\lambda_1} \frac{\lambda_1^x}{x!} \cdot e^{-\lambda_2} \frac{\lambda_2^y}{y!} \quad \text{for } x, y \in \{0, 1, 2, \dots\}$$

(a) As computed above the distribution of Z is given by the convolution. For any $n = 0, 1, 2, \dots$ we have

$$\begin{aligned} P(Z = n) &= P(X + Y = n) \\ &= \sum_{j=0}^n P(X = j) \cdot P(Y = n - j) \\ &= \sum_{j=0}^n e^{-\lambda_1} \frac{\lambda_1^j}{j!} \cdot e^{-\lambda_2} \frac{\lambda_2^{n-j}}{(n-j)!} \\ &= e^{-(\lambda_1 + \lambda_2)} \sum_{j=0}^n \frac{\lambda_1^j \lambda_2^{n-j}}{j!(n-j)!} \\ &= e^{-(\lambda_1 + \lambda_2)} \frac{1}{n!} \sum_{j=0}^n \frac{n!}{j!(n-j)!} \lambda_1^j \lambda_2^{n-j} \\ &= e^{-(\lambda_1 + \lambda_2)} \frac{(\lambda_1 + \lambda_2)^n}{n!} \end{aligned}$$

where in the last line we have used the binomial expansion (2.1.1). Hence we can conclude that $Z \sim \text{Poisson}(\lambda_1 + \lambda_2)$. The above calculation is easily extended by an induction argument to obtain the fact that if $\lambda_i > 0$, X_i , $1 \leq i \leq k$ are independent $\text{Poisson}(\lambda_i)$ distributed random variables (respectively). Then $Z = \sum_{i=1}^k X_i$ has $\text{Poisson}(\sum_{i=1}^k \lambda_i)$ distribution. Thus if we have k independent Poisson (λ) random variables then $\sum_{i=1}^k X_i$ has $\text{Poisson}(k\lambda)$ distribution.

(b) We readily observe that X and Z are dependent. We shall now try to understand the conditional distribution of $(X | Z = n)$. Since the range of X and Y does not have any negative

numbers, given that $Z = X + Y = n$, X can only take values in $\{0, 1, 2, 3, \dots, n\}$. For $k \in \{0, 1, 2, 3, \dots, n\}$ we have

$$\begin{aligned} P(X = k | Z = n) &= \frac{P(X = k, X + Y = n)}{P(X + Y = n)} = \frac{P(X = k, Y = n - k)}{P(X + Y = n)} \\ &= \frac{P(X = k)P(Y = n - k)}{P(X + Y = n)} \\ &= \frac{e^{-\lambda_1} \frac{\lambda_1^k}{k!} \cdot e^{-\lambda_2} \frac{\lambda_2^{n-k}}{(n-k)!}}{e^{-(\lambda_1+\lambda_2)} \frac{(\lambda_1+\lambda_2)^n}{n!}} \\ &= \frac{n!}{k!(n-k)!} \frac{\lambda_1^k \lambda_2^{n-k}}{(\lambda_1 + \lambda_2)^n} \\ &= \binom{n}{k} \left(\frac{\lambda_1}{\lambda_1 + \lambda_2} \right)^k \left(\frac{\lambda_2}{\lambda_1 + \lambda_2} \right)^{n-k}. \end{aligned}$$

Hence $(X | Z = n) \sim \text{Binomial}(n, \frac{\lambda_1}{\lambda_1 + \lambda_2})$. ■

The point of the examples above is that a probability associated with a functional value $f(X_1, X_2, \dots, X_n)$ may be calculated directly from the probabilities associated with the input variables X_1, X_2, \dots, X_n . The following theorem explains how this may be accomplished generally for any number of variables.

THEOREM 3.3.5. *Let X_1, X_2, \dots, X_n be random variables defined on a single sample space S . Let f be a function of n variables for which $f(X_1, X_2, \dots, X_n)$ is defined in the range of the X_j variables. Let B be a subset of the range of f . Then,*

$$P(f(X_1, X_2, \dots, X_n) \in B) = P((X_1, X_2, \dots, X_n) \in f^{-1}(B))$$

Proof - First note that both of the events $(f(X_1, X_2, \dots, X_n) \in B)$ and $((X_1, X_2, \dots, X_n) \in f^{-1}(B))$ are subsets of S since outcomes $s \in S$ determine the values of the X_j variables which in turn determine the output of f . The theorem follows immediately from the set theoretic fact that

$$f(X_1(s), X_2(s), \dots, X_n(s)) \in B \iff (X_1(s), X_2(s), \dots, X_n(s)) \in f^{-1}(B)$$

This is because the expression $f(X_1(s), X_2(s), \dots, X_n(s)) \in B$ is what defines s to be an outcome in the event $(f(X_1, X_2, \dots, X_n) \in B)$. Likewise the expression $(X_1(s), X_2(s), \dots, X_n(s)) \in f^{-1}(B)$ defines s to be in the event $((X_1, X_2, \dots, X_n) \in f^{-1}(B))$. Since these events are equal, they have the same probability. ■

3.3.2 Functions and Independence

If X and Y are independent random variables, does that guarantee that functions $f(X)$ and $g(Y)$ of these random variables are also independent? If we take the intuitive view of independence as saying “knowing information about X does not affect the probabilities associated with Y ” then it seems the answer should be “yes”. After all, X determines the value of $f(X)$ and Y determines the value of $g(Y)$. So information about $f(X)$ should translate to information about X and information about $g(Y)$ should translate to information about Y . Therefore if information about

$f(X)$ affected probabilities associated with $g(Y)$, then it seems there should be information about X that would affect the probability associated with Y . Below we generalize this argument and make it more rigorous.

THEOREM 3.3.6. *Let $n > 1$ be a positive integer. For each $j \in \{1, 2, \dots, n\}$ define a positive integer m_j and suppose $X_{i,j}$ is an array of mutually independent random variables for $j \in \{1, 2, \dots, n\}$ and $i \in \{1, 2, \dots, m_j\}$. Let f_j be functions such that the quantity*

$$Y_j = f_j(X_{1,j}, X_{2,j}, \dots, X_{m_j,j})$$

is defined for the outputs of the $X_{i,j}$ variables. Then the resulting variables Y_1, Y_2, \dots, Y_n are mutually independent.

Informally this theorem says that random quantities produced from independent inputs will, themselves, be independent.

Proof - Let B_1, B_2, \dots, B_n be sets in the ranges of Y_1, Y_2, \dots, Y_n respectively. Use of independence and set-theoretic identities then shows

$$\begin{aligned} P(Y_1 \in B_1, \dots, Y_n \in B_n) &= P(f_1(X_{1,1}, \dots, X_{m_1,1}) \in B_1, \dots, f_n(X_{1,n}, \dots, X_{m_n,n}) \in B_n) \\ &= P((X_{1,1}, \dots, X_{m_1,1}) \in f_1^{-1}(B_1), \dots, (X_{1,n}, \dots, X_{m_n,n}) \in f_n^{-1}(B_n)) \\ &= \prod_{i=1}^n P((X_{i,1}, \dots, X_{m_i,i}) \in f_i^{-1}(B_i)) \\ &= \prod_{i=1}^n P(f_i(X_{i,1}, \dots, X_{m_i,i}) \in B_i) \\ &= P(Y_1 \in B_1) \cdots P(Y_n \in B_n) \end{aligned}$$

which proves that Y_1, Y_2, \dots, Y_n are mutually independent. ■

EXERCISES

Ex. 3.3.1. Let $X \sim \text{Uniform}(\{1, 2, 3\})$ and $Y \sim \text{Uniform}(\{1, 2, 3\})$ be independent and let $Z = X + Y$.

- (a) Determine the range of Z .
- (b) Determine the distribution of Z .
- (c) Is Z uniformly distributed over its range?

Ex. 3.3.2. Consider the experiment of rolling three dice and calculating the sum of the rolls. Answer the following questions.

- (a) What is the range of possible results of this experiment?
- (b) Calculate the probability the sum equals three.
- (c) Calculate the probability the sum equals four.
- (d) Calculate the probability the sum equals five.
- (e) Calculate the probability the sum equals ten.

Ex. 3.3.3. Let $X \sim \text{Bernoulli}(p)$ and $Y \sim \text{Bernoulli}(q)$ be independent.

- (a) Prove that XY is a Bernoulli random variable. What is its parameter?
- (b) Prove that $(1 - X)$ is a Bernoulli random variable. What is its parameter?
- (c) Prove that $X + Y - XY$ is a Bernoulli random variable. What is its parameter?

Ex. 3.3.4. Let $X \sim \text{Binomial}(n, p)$ and $Y \sim \text{Binomial}(m, p)$. Assume X and Y are independent and let $Z = X + Y$. Prove that $Z \sim \text{Binomial}(m + n, p)$.

Ex. 3.3.5. Let $X \sim \text{Negative Binomial}(r, p)$ and $Y \sim \text{Negative Binomial}(s, p)$. Assume X and Y are independent and let $Z = X + Y$. Prove that $Z \sim \text{Negative Binomial}(r + s, p)$.

Ex. 3.3.6. Consider one flip of a single fair coin. Let X denote the number of heads on the flip and let Y denote the number of tails on the flip.

- (a) Show that $X, Y \sim \text{Bernoulli}(\frac{1}{2})$.
- (b) Let $Z = X + Y$ and explain why $P(Z = 1) = 1$.
- (c) Since (b) clearly says that Z cannot be a $\text{Binomial}(2, \frac{1}{2})$, explain why this result does not conflict with the conclusion of Example 3.3.3.

Ex. 3.3.7. Let $X \sim \text{Geometric}(p)$ and $Y \sim \text{Geometric}(p)$ be independent and let $Z = X + Y$.

- (a) Determine the range of Z .
- (b) Use a convolution to prove that $P(Z = n) = (n - 1)p^2(1 - p)^{n-2}$.
- (c) Recall from the discussion of Geometric distributions that $(X = 1)$ is the most likely result for X and $(Y = 1)$ is the most likely result for Y . This does not imply that $(Z = 2)$ is the most likely outcome for Z . Determine the values of p for which $P(Z = 3)$ is larger than $P(Z = 2)$.

Ex. 3.3.8. Let X_1, X_2, X_3, X_4 be an i.i.d. sequence of $\text{Bernoulli}(p)$ random variables. Let $Y = X_1 + X_2 + X_3 + X_4$. Prove that $P(Y = 2) = 6p^2(1 - p)^2$.

Ex. 3.3.9. Let X_1, X_2, \dots, X_n be an i.i.d. sequence of $\text{Bernoulli}(p)$ random variables. Let $Y = X_1 + X_2 + \dots + X_n$. Prove that $Y \sim \text{Binomial}(n, p)$.

Ex. 3.3.10. Let X_1, X_2, \dots, X_r be an i.i.d. sequence of $\text{Geometric}(p)$ random variables. Let $Y = X_1 + X_2 + \dots + X_r$. Prove that $Y \sim \text{Negative binomial}(r, p)$.

Ex. 3.3.11. Let X_1, X_2, X_3, X_4 be an i.i.d. sequence of $\text{Bernoulli}(p)$ random variables. Let $Y = X_1 + X_2$ and let $Z = X_3 + X_4$. Note that Example 3.3.3 guarantees that $Y, Z \sim \text{Binomial}(2, p)$.

- (a) Create a chart describing the joint distribution of Y and Z .
- (b) Use the chart from (a) to explain why Y and Z are independent.
- (c) Explain how you could use Theorem 3.3.6 to reach the conclusion that Y and Z are independent without calculating their joint distribution.

Ex. 3.3.12. Let X_1, X_2, X_3 be an i.i.d. sequence of $\text{Bernoulli}(p)$ random variables. Let $Y = X_1 + X_2$ and let $Z = X_2 + X_3$. Note that Example 3.3.3 guarantees that $Y, Z \sim \text{Binomial}(2, p)$.

- (a) Create a chart describing the joint distribution of Y and Z .

- (b) Use the chart from (a) to explain why Y and Z are not independent.
- (c) Explain why the conclusion from (b) is not inconsistent with Theorem 3.3.6.

Ex. 3.3.13. Let X_1, X_2, \dots, X_n be an i.i.d. sequence of discrete random variables and let Z be the maximum of these n variables. Let r be a real number and let $R = P(X_1 \leq r)$. Prove that $P(Z \leq r) = R^n$.

Ex. 3.3.14. Let X_1, X_2, \dots, X_n be an i.i.d. sequence of discrete random variables and let Z be the minimum of these n variables. Let r be a real number and let $R = P(X_1 \leq r)$. Prove that $P(Z \leq r) = 1 - (1 - R)^n$.

Ex. 3.3.15. Let $X \sim \text{Geometric}(p)$ and let $Y \sim \text{Geometric}(q)$ be independent random variables. Let Z be the smaller of X and Y . It is a fact that Z is also geometrically distributed. This problem asks you to prove this fact using two different methods.

METHOD I:

- (a) Explain why the event $(Z = n)$ can be written as the disjoint union

$$(Z = n) = (X = n, Y = n) \cup (X = n, Y > n) \cup (X > n, Y = n)$$

- (b) Recall from the proof of the memoryless property of geometric random variables that $P(X > m) = \frac{1}{2^m}$. Use this fact and part (a) to prove that

$$P(Z = n) = [(1 - p)(1 - q)]^{n-1} (pq + p(1 - q) + (1 - p)q)$$

- (c) Use (b) to conclude that $Z \sim \text{Geometric}(r)$ for some quantity r and calculate the value of r in terms of the p and q .

METHOD II: Recall that geometric random variables first arose from noting the time it takes for a sequence of Bernoulli trials to first produce a success. With that in mind, let A_1, A_2, \dots be Bernoulli(p) random variables and let B_1, B_2, \dots be Bernoulli(q) random variables. Further assume the A_j and B_k variables collectively are mutually independent. The variable X may be viewed as the number of the first A_j that produces a result of 1 and the variable Y may be viewed similarly for the B_k sequence.

- (a) Let C_j be a random variable that is 1 if either $A_j = 1$ or $B_j = 1$ (or both), and is equal to 0 otherwise. Prove that $C_j \sim \text{Bernoulli}(r)$ for some quantity r and calculate the value of r in terms of p and q .
- (b) Explain why the sequence C_1, C_2, \dots are mutually independent random variables.
- (c) Let Z be the random variable that equals the number of the first C_j that results in a 1 and explain why Z is the smaller of X and Y .
- (d) Use (c) to conclude that $Z \sim \text{Geometric}(r)$ for the value of r calculated in part (a).

Ex. 3.3.16. Each day during the hatching season along the Odisha and Northern Tamil Nadu coast line a Poisson (λ) number of turtle eggs hatch giving birth to young turtles. As these turtles swim into the sea the probability that they will survive each day is p . Assume that number of hatchings on each day and the life of the turtles born are all independent. Let $X_1 = 0$ and for $i \geq 2$, X_i be the total number of turtles alive at sea on the i^{th} morning of the hatching season before the hatchings on the i -th day. Find the distribution of X_n .

5

CONTINUOUS PROBABILITIES AND RANDOM VARIABLES

We have thus far restricted discussion to discrete spaces and discrete random variables – those consisting of at most a countably infinite number of outcomes. This is not because it is not possible, interesting, or useful to consider probabilities on an uncountably infinite set such as the real line or the interval $(0, 1)$. Instead, there are a few technicalities that arise when discussing such probabilities that are best avoided until they are needed. That time is now.

5.1 UNCOUNTABLE SAMPLE SPACES AND DENSITIES

Suppose we want to randomly select a number on the interval $(0, 1)$ in some uniform way. In the discrete setting we would have said that “uniform” meant that every outcome in our sample space $S = (0, 1)$ was equally likely. Suppose we took that same approach here and declared that there was some value p for which $P(\{x\}) = p$ for every $x \in (0, 1)$. Then if we let E be the event $E = \{\frac{1}{n} : n = 2, 3, 4, \dots\} \subset S$, we find that

$$\begin{aligned} P(E) &= P(\{\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots\}) \\ &= P(\frac{1}{2}) + P(\frac{1}{3}) + P(\frac{1}{4}) + \dots \\ &= p + p + p + \dots \end{aligned}$$

If $p > 0$ this sum diverges to infinity, which cannot be since it describes a probability. Therefore it must be that $p = 0$. If every individual outcome in $S = (0, 1)$ is equally likely, then each outcome must have a probability of zero. After several chapters considering only discrete probabilities many readers may suspect that this, in and of itself, is a contradiction. How is it possible for $P(S) = 1$ when every single element of S has probability zero? Could not one then show

$$\begin{aligned} P(S) &= P(\bigcup_{s \in S} \{s\}) \\ &= \sum_{s \in S} P(\{s\}) \\ &= \sum_{s \in S} 0 \\ &= 0 \end{aligned}$$

using the probability axioms? The answer to that question is “no”. The probability space axiom that allows us to write the probability of a disjoint union as the sum of separate probabilities only applies to countable collections of events. But the events $\{s\}$ that combine to create $(0, 1)$ are an uncountable collection. If S is uncountable, we could still have $P(S) = 1$ even if every individual element of $s \in S$ has probability zero.

However, all of that does not yet explain how to define a uniform probability on $(0, 1)$. Knowing that each individual outcome has probability zero does not tell us how to calculate $P([\frac{1}{4}, \frac{3}{4}])$, for example, since we cannot simply add up the probabilities of each of the constituent outcomes individually. Instead we need to reinterpret what we mean by “uniform” in this situation. It would make sense to suggest that the event $[\frac{1}{4}, \frac{3}{4}]$ should have a probability of $\frac{1}{2}$ since its length is exactly

half of the length of $(0, 1)$. Indeed it is tempting (and essentially correct) to declare that $P(A)$ should be the length of the set A . The complication with making such a statement is that, although length is easy to define if A is an interval or even a countable collection of disjoint intervals, it is not even possible to consistently define a length for every single subset of $(0, 1)$. Because of this unfortunate fact, we will need to reconsider which subsets of S are actually events which will be assigned a probability.

At a minimum we will want events to include any interval. The axioms and basic properties of probability spaces also require that for any collection of events we must be able to consider complements and countable unions of these events. Further, the entire sample space S should also be considered a legitimate event. Consequently we make the following definition.

DEFINITION 5.1.1. (σ -field) *If S is a sample space, then a σ -field \mathcal{F} is a collection of subsets of S such that*

$$(1) \quad S \in \mathcal{F}$$

$$(2) \quad \text{If } A \in \mathcal{F} \text{ then } A^c \in \mathcal{F}$$

$$(3) \quad \text{If } A_1, A_2, \dots \text{ is a countable collection of sets in } \mathcal{F} \text{ then } \bigcup_{n=1}^{\infty} A_n \in \mathcal{F}$$

We shall refer to an element of the σ -field as an event.

If S happens to be the set of real numbers there is a smallest σ -field that contains all intervals, and this collection of subsets of \mathbb{R} is known as the Borel sets. It happens that the concept of the “length” of a set can be consistently described for such sets. Because of this we will modify our definition of probability space slightly at this point.

DEFINITION 5.1.2. (Probability Space Axioms) *Let S be a sample space and let \mathcal{F} be a σ -field of S . A “probability” is a function $P : \mathcal{F} \rightarrow [0, 1]$ such that*

$$(1) \quad P(S) = 1;$$

$$(2) \quad \text{If } E_1, E_2, \dots \text{ are a countable collection of disjoint events in } \mathcal{F}, \text{ then}$$

$$P\left(\bigcup_{j=1}^{\infty} E_j\right) = \sum_{j=1}^{\infty} P(E_j).$$

The triplet (S, \mathcal{F}, P) is referred to as a probability space.

Our old definition is simply a special case where the σ -field was the collection of all subsets of S , so all results we have previously seen in the discrete setting are still legitimate in this new framework. There are many technicalities that arise due to the fact that not every set may be viewed as an event, but these issues would be distracting from the primary goal of this text. Thus we give the definitions above only to provide the modern definition of probability space.

Throughout the remainder of the sections on continuous probability spaces we will restrict our attention to the sample space being \mathbb{R} . Whenever we state or prove anything for an event A (a

Borel set) we shall restrict ourselves to the case the event A is a finite or countable unions of intervals. This will enable us to use standard results from calculus and thereby avoid technicalities. A thorough study of the Borel sets and the related theory of integration is beyond the scope of this text (the interested reader may see [AS09] in the bibliography for additional information).

5.1.1 Probability Densities on \mathbb{R}

The primary way we will define continuous probabilities on \mathbb{R} is through a “density function”. We begin by providing an example of what should be meant by a uniform distribution on $(0, 1)$.

EXAMPLE 5.1.3. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a function defined by

$$f(x) = \begin{cases} 1 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

For an event A define $P(A) = \int_A f(x) dx$. Note that for an interval $A = [a, b] \subset (0, 1)$ it happens that $P(A)$ is just the length of the interval.

$$\begin{aligned} P(A) &= \int_A f(x) dx \\ &= \int_a^b 1 dx \\ &= b - a \end{aligned}$$

For disjoint unions of intervals, the lengths simply add. For instance if $A = [\frac{1}{5}, \frac{2}{5}] \cup [\frac{3}{5}, \frac{4}{5}]$, then

$$\begin{aligned} P(A) &= \int_{[\frac{1}{5}, \frac{2}{5}] \cup [\frac{3}{5}, \frac{4}{5}]} f(x) dx \\ &= \int_{\frac{1}{5}}^{\frac{2}{5}} 1 dx + \int_{\frac{3}{5}}^{\frac{4}{5}} 1 dx \\ &= \frac{1}{5} + \frac{1}{5} = \frac{2}{5} \end{aligned}$$

which is the sum of the lengths of the two component intervals. In particular note that $P((0, 1)) = 1$ while $P(\{c\}) = 0$ for any c since a single point has no length. Similarly, if $A = [a, b]$ is an interval that is disjoint from $(0, 1)$, then

$$\begin{aligned} P(A) &= \int_A f(x) dx \\ &= \int_a^b 0 dx \\ &= 0 \end{aligned}$$

We will soon see that P defines a probability on \mathbb{R} . From the computation above this probability gives equal likelihood to all equal-width intervals within $(0, 1)$ and assigns zero probability to any interval outside of $(0, 1)$. Therefore it is consistent with the properties a uniform probability on $(0, 1)$ should have. ■

The function f from the example above is known as a density. What properties must be required of such a function in order for it to define a probability? The fact that probabilities

cannot be negative suggests we will need to require $f(x)$ to be non-negative for all real numbers x . The fact that $P(S) = 1$ means that $\int_{-\infty}^{\infty} f(x) dx$ has to be 1. It turns out that these two requirements are essentially all that are needed. The only other assumption we will make is that a density function be piecewise continuous. Though this final requirement is more restrictive than necessary, the assumption will help avoid technicalities and will include all densities of interest to us in the remainder of the text. We give a precise definition.

DEFINITION 5.1.4. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be called a density function if f satisfies the following:

- (i) $f(x) \geq 0$,
- (ii) f is piecewise-continuous, and
- (iii) $\int_{-\infty}^{\infty} f(x) dx = 1$.

We proceed to state and prove a result that will help us construct probabilities on \mathbb{R} with the help of density functions. This will also ensure that in Example 5.1.3 we indeed constructed a probability on \mathbb{R} .

THEOREM 5.1.5. Let $f(x)$ be a density function. Define

$$P(A) = \int_A f(x) dx, \quad (5.1.1)$$

for any event $A \subset \mathbb{R}$. Then P defines a probability on \mathbb{R} . The function f is called the “density function” for the probability P .

Proof - First note

$$\begin{aligned} P(\mathbb{R}) &= \int_{\mathbb{R}} f(x) dx \\ &= \int_{-\infty}^{\infty} f(x) dx = 1 \end{aligned}$$

by assumption, so the entire sample space has probability 1. Now let A be a Borel subset of \mathbb{R} . Since $f(x)$ is non-negative,

$$\begin{aligned} P(A) &= \int_A f(x) dx \geq 0, \quad \text{and} \\ P(A) &= \int_A f(x) dx \leq \int_{\mathbb{R}} f(x) dx = 1, \end{aligned}$$

so $P(A) \in [0, 1]$. Finally, if E_1, E_2, \dots are a countable collection of disjoint events, then

$$\begin{aligned} P\left(\bigcup_{n=1}^{\infty} E_n\right) &= \int_{\bigcup_{n=1}^{\infty} E_n} f(x) dx \\ &= \sum_{n=1}^{\infty} \int_{E_n} f(x) dx \\ &= \sum_{n=1}^{\infty} P(E_n). \end{aligned}$$

Therefore P satisfies the conditions of a probability space on \mathbb{R} . ■

EXAMPLE 5.1.6. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} 3x^2 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

f is piecewise continuous, $f(x)$ is non-negative for all x and

$$\int_{\mathbb{R}} f(x) dx = \int_0^1 3x^2 dx = x^3 \Big|_0^1 = 1.$$

As it satisfies (i) – (iii) in Definition 5.1.4, f is a density function. Let P be as defined in (5.1.1). As with the uniform example, $f(x)$ is zero outside of $(0, 1)$, so events lying outside this interval will have zero probability. However, note that

$$\begin{aligned} P\left(\left[\frac{1}{5}, \frac{2}{5}\right]\right) &= \int_{\frac{1}{5}}^{\frac{2}{5}} 3x^2 dx = \frac{7}{125} \quad \text{while} \\ P\left(\left[\frac{3}{5}, \frac{4}{5}\right]\right) &= \int_{\frac{3}{5}}^{\frac{4}{5}} 3x^2 dx = \frac{37}{125}. \end{aligned}$$

In other words, intervals of the same length do not have equal probabilities; this probability is not uniform on $(0, 1)$.

The probability of individual points is still zero, so $P(\{\frac{1}{5}\}) = P(\{\frac{2}{5}\}) = 0$, but in terms of the density function, $f(\frac{2}{5})$ is four times as large as $f(\frac{1}{5})$. What does this mean in practical terms?

Let ϵ be a small positive quantity (certainly less than $\frac{1}{5}$). Then

$$\begin{aligned} P\left(\left[\frac{1}{5} - \epsilon, \frac{1}{5} + \epsilon\right]\right) &= \frac{2}{25}\epsilon + 2\epsilon^3 \approx \frac{2}{25}\epsilon \quad \text{while} \\ P\left(\left[\frac{2}{5} - \epsilon, \frac{2}{5} + \epsilon\right]\right) &= \frac{8}{25}\epsilon + 2\epsilon^3 \approx \frac{8}{25}\epsilon. \end{aligned}$$

The fact that $f(\frac{2}{5})$ is four times as large as $f(\frac{1}{5})$ essentially means that a tiny interval around $\frac{2}{5}$ has approximately four times the probability of a similarly sized interval around $\frac{1}{5}$. ■

EXERCISES

Ex. 5.1.1. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} 2x & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

(a) Show that f is a probability density function.

(b) Use f to calculate $P((0, \frac{1}{2}))$.

Ex. 5.1.2. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

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

(a) Sketch a graph of the function f .

- (b) Show that f is a probability density function.
 (c) Use f to calculate : $P((0, \frac{1}{4}), P((\frac{3}{2}, 2)), P((-3, -2))$ and $P((\frac{1}{2}, \frac{3}{2}))$.

Ex. 5.1.3. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} k & \text{if } 0 < x < \frac{1}{4} \\ 2k & \text{if } \frac{1}{4} \leq x < \frac{3}{4} \\ 3k & \text{if } \frac{3}{4} \leq x < 1 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find k that makes f a probability density function.
 (b) Sketch a graph of the function f .
 (c) Use f to calculate : $P((0, \frac{1}{4}), P((\frac{1}{4}, \frac{3}{4})), P((\frac{3}{4}, 1))$.

Ex. 5.1.4. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} k \cdot \sin(x) & \text{if } 0 < x < \pi \\ 0 & \text{otherwise} \end{cases}$$

- (a) Determine the value of k that makes f a probability density function.
 (b) Calculate $P((0, \frac{1}{2}))$ and $P((\frac{1}{2}, 1))$.
 (c) Which will be larger, $P((0, \frac{1}{4}))$ or $P((\frac{1}{4}, \frac{1}{2}))$? Explain how you can answer this question without actually calculating either probability.
 (d) A game is played in the following way. A random variable X is selected with a density described by f above. You must select a number r and you win the game if the random variable results in an outcome in the interval $(r - 0.01, r + 0.01)$. Explain how you should choose r to maximize your chance of winning the game. (A formal proof requires only basic calculus, but it should take very little computation to determine the correct answer).

Ex. 5.1.5. Let $\lambda > 0$ and $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } 0 < x \\ 0 & \text{otherwise} \end{cases}$$

- (a) Show that f is a probability density function.
 (b) Let $a > 0$. Find $P((a, \infty))$.

Ex. 5.1.6. Let $a, b \in \mathbb{R}$ and $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$

- (a) Show that f is a probability density function.
 (b) Show that if $I, J \subset [a, b]$ are two intervals that have the same length, then $P(I) = P(J)$.

Ex. 5.1.7. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} \frac{1}{x^2} & \text{if } 1 < x \\ 0 & \text{otherwise} \end{cases}$$

- (a) Show that f is a probability density function.
 (b) Let $a > 1$. Find $P((a, \infty))$.

Ex. 5.1.8. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \begin{cases} \frac{1}{6}x^2e^{-x} & \text{if } 0 < x \\ 0 & \text{otherwise} \end{cases}$$

Show that f is a probability density function.

Ex. 5.1.9. For any $x \in \mathbb{R}$, the hyperbolic secant is defined as $\operatorname{sech} x = \frac{2}{(e^x + e^{-x})}$. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(x) = \frac{1}{2} \operatorname{sech}\left(\frac{\pi}{2}x\right), x \in \mathbb{R}$$

Show that f is a probability density function.

Ex. 5.1.10. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be defined by

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

Follow the steps below to show that the function f is a density function.

- (a) Let $I = \int_{-\infty}^{\infty} e^{-x^2/2} dx$ and then explain why

$$I^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)/2} dx dy$$

(Hint: Write I^2 as a product of two integrals each over a different variable and explain why the resulting expression may be written as the double integral above).

- (b) Explain why

$$I^2 = \int_0^{2\pi} \int_0^{\infty} r \cdot e^{-r^2/2} dr d\theta$$

after switching from rectangular to polar coordinates. (Hint: Use the fact from multivariate calculus that after the change of variables $(dx dy)$ becomes $(r dr d\theta)$ and explain the new limits of integration based on the region being described in the plane).

- (c) Compute the integral from (b) to find the value of I .

- (d) Use (c) to show that $\int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1$. (Hint: Use a u-substitution $u = \frac{x-\mu}{\sigma}$).

5.2 CONTINUOUS RANDOM VARIABLES

Just as the move from discrete to continuous spaces required a slight change in the definition of probability space, so it also requires a slight change in the definition of random variable. In the discrete setting we frequently needed to consider the preimage $X^{-1}(A)$ of a set. Now we need to make sure that such a preimage is a legitimate event.

DEFINITION 5.2.1. Let (S, \mathcal{F}, P) be a probability space and let $X : S \rightarrow \mathbb{R}$ be a function. Then X is a random variable provided that whenever B is an event in \mathbb{R} (i.e. a Borel set), $X^{-1}(B)$ is also an event in \mathcal{F} .

Note that in the discrete setting this extra condition was met trivially as every subset of S was an event. Therefore the discrete setting is simply a special case of this new definition. As with the introduction of σ -fields, we include this definition for completeness. We will only consider functions which meet this criterion. In this section we shall consider only continuous random variables. These are defined next.

DEFINITION 5.2.2. Let (S, \mathcal{F}, P) be a probability space. A random variable $X : S \rightarrow \mathbb{R}$ is called a continuous random variable if there exists a density function $f_X : \mathbb{R} \rightarrow \mathbb{R}$ such that for any event A in \mathbb{R} ,

$$P(X \in A) = \int_A f_X(x) dx. \quad (5.2.1)$$

f_X is called the probability density function of X .

The following lemma demonstrates an elementary property of continuous random variables that distinguishes them from discrete random variables.

LEMMA 5.2.3. Let X be a continuous random variable. Then for any $a \in \mathbb{R}$,

$$P(X = a) = 0 \quad (5.2.2)$$

Proof- Let $a \in \mathbb{R}$, then $P(X = a) = \int_a^a f(x) dx = 0$. ■

Random variables may also be described using a “distribution function” (also commonly known as a “cumulative distribution function”).

DEFINITION 5.2.4. If X is a random variable then its distribution function $F : \mathbb{R} \rightarrow [0, 1]$ is defined by

$$F(x) = P(X \leq x). \quad (5.2.3)$$

When it must be emphasized that a distribution function belongs to a particular random variable X the notation $F_X(x)$ will be used to indicate the random variable.

Though a distribution function is defined for any real-valued random variable, there is a special relationship between $f_X(x)$ and $F_X(x)$ when the random variable has a density.

THEOREM 5.2.5. Let X be a random variable with a piecewise continuous density function $f(x)$. If $F(x)$ denotes the distribution function of X then

$$F(x) = \int_{-\infty}^x f(x) dx. \quad (5.2.4)$$

Moreover, where $f(x)$ is continuous, $F(x)$ is differentiable and $F'(x) = f(x)$.

Proof - By definition $F(x) = P(X \leq x) = P(X \in (-\infty, x])$, but this probability is described in terms of an integral over the density of X , so $F(x) = \int_{-\infty}^x f(x) dx$.

The result that $F'(x) = f(x)$ then follows from the fundamental theorem of calculus after taking derivatives of both sides of the equation (when such a derivative exists). Note, in particular, that since densities are assumed to be piecewise continuous, their corresponding distribution functions are piecewise differentiable. \blacksquare

This theorem will be useful for computation, but it also shows that the distribution of a continuous random variable X is completely determined by its distribution function F_X . That is, if we know $F_X(x)$ and want to calculate $P(X \in A)$ for some set A we could do so by differentiating $F_X(x)$ to find $f_X(x)$ and then integrating this density over the set A . In fact $F_X(x)$ always completely determines the distribution of X (regardless of whether or not X is a continuous random variable), but a proof of that fact is beyond the scope of the course and will not be needed for subsequent results.

5.2.1 Common Distributions

In the literature random variables whose distributions satisfy (5.2.1) are called absolutely continuous random variables and those that satisfy (5.2.2) are referred to as continuous random variables. Since we shall only consider continuous random variables that satisfy (5.2.1) we refer to them as continuous random variables.

There are many continuous distributions that commonly arise. Some of these are continuous analogs of discrete random variables we have already studied. We will define these in the context of continuous random variables having the corresponding distributions. We begin with the already discussed uniform distribution but on an arbitrary interval.

DEFINITION 5.2.6. $X \sim \text{Uniform}(a, b)$: Let (a, b) be an open interval. If X is a random variable with its probability density function given by

$$f(x) = \begin{cases} \frac{1}{(b-a)} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$

then X is said to be uniformly distributed on (a, b) . Note that this is consistent with the example at the beginning of the section since the density of a $\text{Uniform}(0, 1)$ is one on the interval $(0, 1)$ and zero elsewhere. Further, recall that in Exercise 5.1.6 we have shown that f is indeed a probability density function.

Since X only takes values on (a, b) if $x < a$ then $P(X \leq x) = 0$ while if $x > b$ then $P(X \leq x) = 1$. So let $a \leq x \leq b$. Then,

$$P(X \leq x) = \int_{-\infty}^x f_X(y) dy = \int_{-\infty}^a 0 dy + \int_a^x \frac{1}{b-a} dy = \frac{x-a}{b-a}.$$

Therefore the distribution function for X is

$$F_X(x) = \begin{cases} 0 & \text{if } x < a \\ \frac{x-a}{b-a} & \text{if } a \leq x \leq b \\ 1 & \text{if } x > b \end{cases}$$

Exponential Random Variable

The next continuous distribution we introduce is called the exponential distribution. It is well known from physical experiments that the radioactive isotopes decay to its stable form. Suppose there were $N(0)$ atoms of a certain radioactive material at time 0 then one is interested in the fraction of radioactive material that have not decayed at time $t > 0$. It is observed from experiments that if $N(t)$ is the number of atoms of radioactive material that has not decayed by time t then the fraction

$$\frac{N(t)}{N(0)} \approx e^{-\lambda t},$$

for some $\lambda > 0$. One can introduce a probability model for the above experiment in the following manner. Suppose X represented the time taken by a randomly chosen radioactive atom to decay to its stable form. The distribution of the random variable X needs to satisfy

$$P(X \geq t) = e^{-\lambda t}, \quad (5.2.5)$$

for $t > 0$. It is possible to define such a random variable.

DEFINITION 5.2.7. $X \sim \text{Exp}(\lambda)$: Suppose $\lambda > 0$. If X is a random variable with its probability density function given by

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

it is said to be distributed exponentially with parameter λ . Recall that in Exercise 5.1.5 we have shown that f is indeed a probability density function. Since X only takes values on $(0, \infty)$ if $x < 0$ then $P(X \leq x) = 0$. So let $x \geq 0$. Then,

$$P(X \leq x) = \int_{-\infty}^x f_X(x) dx = \int_{-\infty}^0 0 dx + \int_0^x \lambda e^{-\lambda y} dy = -e^{-\lambda y} \Big|_0^x = 1 - e^{-\lambda x}.$$

Therefore the distribution function for X is

$$F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 - e^{-\lambda x} & \text{if } 0 \leq x \end{cases}$$

We have previously seen that geometric random variables have the memoryless property (See (3.2.2)). It turns out that the exponential random variable also possess the *memoryless property* in continuous time. Clearly if $X \sim \text{Exp}(\lambda)$ then $P(X \geq 0) = 1$ and

$$P(X \geq t) = P(X \in [t, \infty)) = \int_t^\infty \lambda e^{-\lambda x} dx = -e^{-\lambda x} \Big|_t^\infty = e^{-\lambda t},$$

for $t > 0$. Further if $s, t > 0$, $X > s + t$ imples $X > s$. So

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

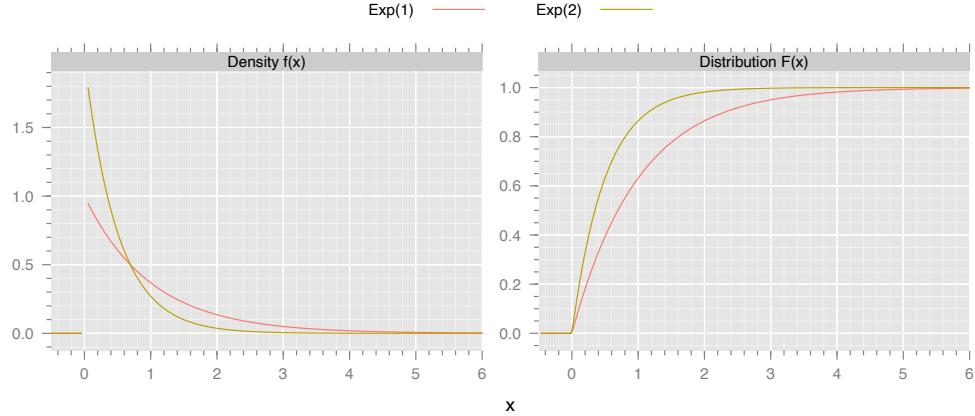


Figure 5.1: The shape of typical Exponential density and cumulative distribution functions.

Therefore for all $s, t > 0$

$$P(X > s + t | X > s) = P(X > t) \quad (5.2.6)$$

Thinking of the variables s and t in terms of time, this says that if an exponential random variable has not yet occurred by time s , then its distribution from that time onward continues to be distributed like an exponential random variable with the same parameter. Situations that involve waiting times such as the lifetime of a light bulb or the time spent in a queue at a service counter are often modelled with the exponential distribution. It is a fact (see Exercise 5.2.12) that if a positive continuous random variable has the memoryless property then it necessarily is an exponential random variable.

EXAMPLE 5.2.8. Let $X \sim \text{Exp}(2)$. Calculate the probability that X produces a value larger than 4.

The density of X is

$$f(x) = \begin{cases} 2e^{-2x} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

So, $P(X > 4)$ may be calculated via an integral.

$$\begin{aligned} P(X > 4) &= \int_4^\infty 2e^{-2x} dx \\ &= -e^{-2x} \Big|_4^\infty = 0 - (-e^{-8}) = e^{-8} \approx 0.000335 \end{aligned}$$

So there is only about a 0.0335% chance of such a result. ■

Normal Random Variable

Of all continuous distributions, The normal distribution (also sometimes called a “Gaussian distribution”) is the most fundamental for applications of statistics as it frequently arises as a limiting distribution of sampling procedures.

DEFINITION 5.2.9. $X \sim \text{Normal}(\mu, \sigma^2)$: Let $\mu \in \mathbb{R}$ and let $\sigma > 0$. Then X is said to be normally distributed with parameters μ and σ^2 if it has the density

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (5.2.7)$$

for all $x \in \mathbb{R}$. We will prove that μ and σ are, respectively, the mean and standard deviation of such a random variable (See Definition 6.1.1, Definition 6.1.9, Example 6.1.11). Recall that in Exercise 5.1.10 we have seen that f is a probability density function.

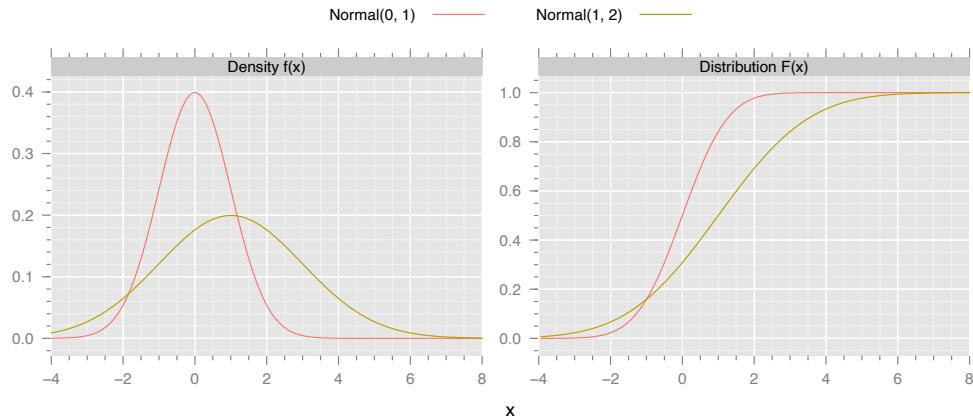


Figure 5.2: The shape of typical Normal density and cumulative distribution functions.

It is observed during statistical experiments that if X were to denote the number of leaves in an apple tree or the height of adult men in a population then X would be close to a normal random variable with appropriate parameters μ and σ^2 . It also arises as a limiting distribution. We shall discuss this aspect in general in Chapter 8, but here we will briefly mention one such limit that appears as an approximation for Binomial probabilities.

Suppose we have X_1, X_2, \dots, X_n are i.i.d Bernoulli (p) random variables. Then we know that $S_n = \sum_{i=1}^n X_i$ is a Binomial (n, p) random variable. In Theorem 2.2.2 we saw that for $\lambda > 0$, $k \geq 1$, $0 \leq p = \frac{\lambda}{n} < 1$,

$$\lim_{n \rightarrow \infty} P(S_n = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

Such an approximation was useful when p was decreasing to zero while n grew to infinity with np remaining constant. The De Moivre-Laplace Central Limit Theorem allows us to consider another form of limit where p remains fixed, but n increases.

THEOREM 5.2.10. (De Moivre-Laplace Central Limit Theorem) Suppose $S_n \sim \text{Binomial}(n, p)$, where $0 < p < 1$. Then for any $a < b$

$$\lim_{n \rightarrow \infty} P(a < \frac{S_n - np}{\sqrt{np(1-p)}} \leq b) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-\frac{x^2}{2}} dx \quad (5.2.8)$$

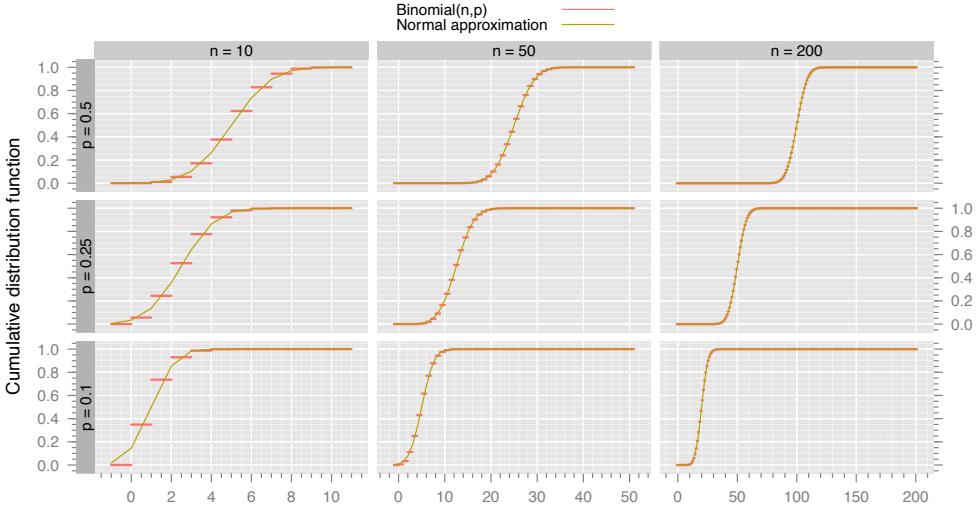


Figure 5.3: The normal approximation to binomial.

We shall omit the proof of the above Theorem for now. We prove it in a more general setting in Chapter 8. For the students well versed with Real Analysis the proof is sketched in Exercise 5.2.16. We refer the reader to [Ram97] for a detailed discussion of the Theorem 5.2.10.

Calculating Normal Probabilities and Necessity of Normal Tables

In a standard introduction to integral calculus one learns many different techniques for calculating integrals. But there are some functions whose indefinite integral has no closed-form solution in terms of simple functions. The density of a normal random variable is one such function. Because of this if $X \sim \text{Normal}(0, 1)$ the probability

$$P(X \leq x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

cannot be expressed exactly in terms of standard functions. Many scientific calculators will have a feature that allows this expression to be evaluated. For example, in R, the command `pnorm(x)` evaluates the integral above. Another common solution in statistical texts is to provide a table of values.

Table 5.1 gives values only for positive values of z because for negative z , $P(X \leq z)$ can be easily computed using the symmetry of the $\text{Normal}(0, 1)$ distribution as (see Figure 5.4)

$$P(X \leq z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \int_{-z}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = 1 - P(X \leq -z). \quad (5.2.9)$$

A more complete version of this table is given in the Appendix. A similar computation can be made for other normally distributed random variables by normalizing them. Suppose $Y \sim \text{Normal}(\mu, \sigma^2)$ and we were interested in finding the distribution function of Y . Observe that

$$P(Y \leq y) = \int_{-\infty}^y \frac{1}{\sigma\sqrt{2\pi}} e^{-(z-\mu)^2/2\sigma^2} dz.$$

	0.00	0.02	0.04	0.06	0.08	0.10	0.12	0.14	0.16	0.18
0.0	0.500	0.508	0.516	0.524	0.532	0.540	0.548	0.556	0.564	0.571
0.2	0.579	0.587	0.595	0.603	0.610	0.618	0.626	0.633	0.641	0.648
0.4	0.655	0.663	0.670	0.677	0.684	0.691	0.698	0.705	0.712	0.719
0.6	0.726	0.732	0.739	0.745	0.752	0.758	0.764	0.770	0.776	0.782
0.8	0.788	0.794	0.800	0.805	0.811	0.816	0.821	0.826	0.831	0.836
1.0	0.841	0.846	0.851	0.855	0.860	0.864	0.869	0.873	0.877	0.881
1.2	0.885	0.889	0.893	0.896	0.900	0.903	0.907	0.910	0.913	0.916
1.4	0.919	0.922	0.925	0.928	0.931	0.933	0.936	0.938	0.941	0.943
1.6	0.945	0.947	0.949	0.952	0.954	0.955	0.957	0.959	0.961	0.962
1.8	0.964	0.966	0.967	0.969	0.970	0.971	0.973	0.974	0.975	0.976
2.0	0.977	0.978	0.979	0.980	0.981	0.982	0.983	0.984	0.985	0.985

Table 5.1: Table of $\text{Normal}(0, 1)$ probabilities. For $X \sim \text{Normal}(0, 1)$, the table gives values of $P(X \leq z)$ for various values of z between 0 and 2.18 upto three digits. The value of z for each entry is obtained by adding the corresponding row and column labels.

Now perform a change of variables $u = \frac{z-\mu}{\sigma}$ so that $du = \frac{1}{\sigma} dz$. This integral then becomes

$$P(Y \leq y) = \int_{-\infty}^{\frac{y-\mu}{\sigma}} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du = P(X \leq \frac{y-\mu}{\sigma}), \quad (5.2.10)$$

where $X \sim \text{Normal}(0, 1)$. Now we may use Table 5.1 to compute the distribution function of Y . We conclude this section with two examples.

EXAMPLE 5.2.11. If $X \sim \text{Normal}(0, 1)$, how likely is it that X will be within one standard deviation of its expected value?

In this case the expected value of the random variable is zero and the standard deviation is one. Therefore the answer is given by

$$\begin{aligned} P(-1 \leq X \leq 1) &= \int_{-1}^1 \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \\ &= \int_{-\infty}^1 \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx - \int_{-\infty}^{-1} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \\ &= P(X \leq 1) - P(X \leq -1) \end{aligned}$$

R tells us that

```
> pnorm(1)
[1] 0.8413447
> pnorm(-1)
[1] 0.1586553
> pnorm(1) - pnorm(-1)
[1] 0.6826895
```

Alternatively, using Table 5.1, we see that $P(X \leq 1) = 0.841$ (upto three decimal places), and by symmetry $P(X \leq -1) = P(X \geq 1) = 1 - P(X \leq 1) = 1 - 0.841 = 0.159$. Therefore, $P(-1 \leq X \leq 1) \approx 0.841 - 0.159 = 0.682$. In other words, there is roughly a 68% chance that a standardized normal random variable will produce a value within one standard deviation of expected value. ■



Figure 5.4: Computation of $\text{Normal}(0, 1)$ probabilities as area under the normal density curve. For $\text{Normal}(0, 1)$ and in fact for any symmetric distribution in general, it is enough to know the distribution function for positive values (see Exercise 5.2.8).

EXAMPLE 5.2.12. A machine fills bags with cashews. The intended weight of cashews in the bag is 200 grams. Assume the machine has a tolerance such that the actual weight of the cashews is a normally distributed random variables with an expected value of 200 grams and a standard deviation of 4 grams. How likely is it that a bag filled by this machine will have fewer than 195 grams of cashews in it?

We know $Y \sim \text{Normal}(200, 4^2)$ and we want the probability $P(Y < 195)$. By above computation, (5.2.10),

$$P(Y < 195) = P\left(X < \frac{195 - 200}{4}\right) = P\left(X < -\frac{5}{4}\right)$$

where $X \sim \text{Normal}(0, 1)$. If we were to use Table 5.1, we would first obtain

$$P\left(X < -\frac{5}{4}\right) = 1 - P\left(X < \frac{5}{4}\right) = 1 - P(X < 1.25) = 1 - 0.896 = 0.104$$

By the R command “pnorm(-5/4)” we obtain 0.1056498. That is, there is slightly more than a 10% chance of a bag this light being produced by the machine. ■

5.2.2 A word about individual outcomes

We began this section by noting that continuous random variables must necessarily give probability zero to any single outcome. It is an awkward consequence of this that two different densities may give rise to exactly the same probabilities. For instance, the functions

$$f(x) = \begin{cases} 1 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

and

$$g(x) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

are different because they assign different values to the points $x = 0$ and $x = 1$. However, these individual points cannot affect the computation of probabilities so both $f(x)$ and $g(x)$ give rise to the same probability distribution. The same thing would occur even if $f(x)$ and $g(x)$ differed in a countably infinite number of points, since these will still have probability zero when taken collectively.

Because of this we will describe $f(x)$ and $g(x)$ as the same density (and sometimes even write $f(x) = g(x)$) when the two densities produce the same probabilities. We do this even when f and g may technically be different functions. Though it is a more restrictive assumption than is necessary, we have required densities to be piecewise continuous. As a consequence of the explanation above, altering the values of the function at the endpoints of intervals of continuity will not change the resulting probabilities and will result in the same density.

EXERCISES

Ex. 5.2.1. Suppose X was continuous random variable with distribution function F . Express the following probabilities in terms of F :

- (a) $P(a < X \leq b)$, where $-\infty < a < b < \infty$
- (b) $P(a < X < \infty)$ where $a \in \mathbb{R}$.
- (c) $P(|X - a| \geq b)$ where $a, b \in \mathbb{R}$ and $b > 0$.

Ex. 5.2.2. Let $R > 0$ and $X \sim \text{Uniform} [0, R]$. Let $Y = \min(X, \frac{R}{10})$. Find the distribution function of Y .

Ex. 5.2.3. Let X be a random variable with distribution function given by

$$F(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } 0 < x < \frac{1}{4} \\ \frac{x}{2} + \frac{1}{8} & \text{if } \frac{1}{4} \leq x < \frac{3}{4} \\ 2x - 1 & \text{if } \frac{3}{4} \leq x < 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

- (a) Sketch a graph of the function F .
- (b) Use F to calculate : $P([0, \frac{1}{4}])$, $P([\frac{1}{8}, \frac{3}{2}])$, $P((\frac{3}{4}, \frac{7}{8}])$.
- (c) Find the probability density function of X .

Ex. 5.2.4. Let X be a continuous random variable with distribution function $F : \mathbb{R} \rightarrow [0, 1]$. Then $G : \mathbb{R} \rightarrow [0, 1]$ given by

$$G(x) = 1 - F(x)$$

is called the reliability function of X or the right tail distribution function of X . Suppose $T \sim \text{Exponential}(\lambda)$ for some $\lambda > 0$, then find the reliability function of T .

Ex. 5.2.5. Let X be a random variable whose probability density function $f : \mathbb{R} \rightarrow [0, 1]$ is given by

$$f(x) = \begin{cases} kx^{k-1}e^{-x^k} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find the distribution function of X for $k = 2$.

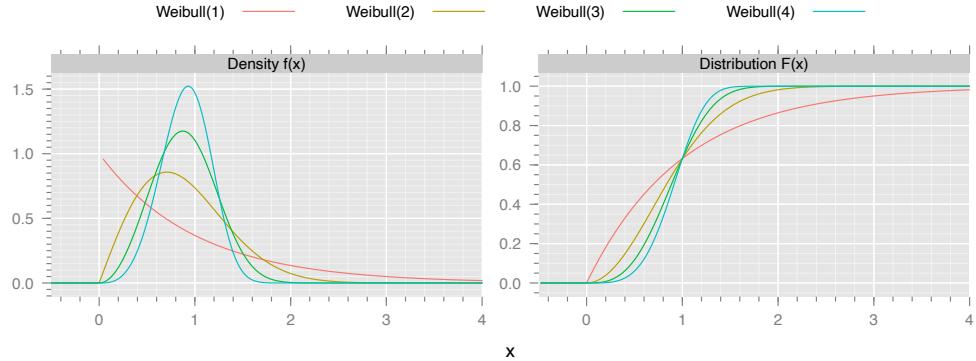


Figure 5.5: The shape of typical Weibull density and cumulative distribution functions.

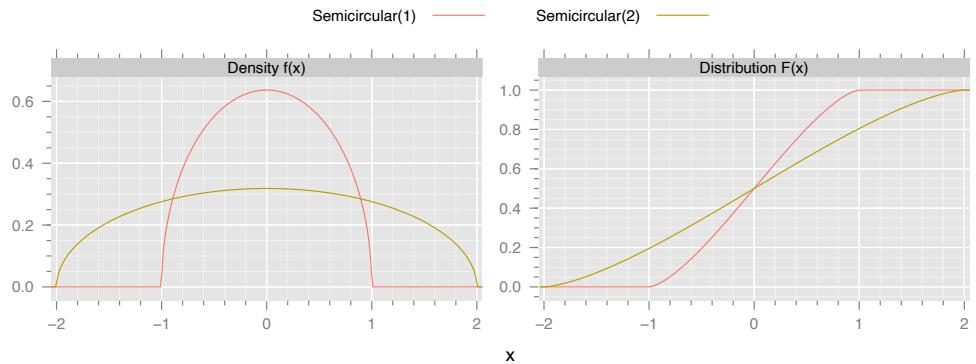


Figure 5.6: The shape of the semicircular density and cumulative distribution functions.

(b) Find the distribution function of X for general k .

The distribution of X is called the Weibull distribution. Figure 5.5 plots the Weibull distribution for selected values of k .

Ex. 5.2.6. Let X be a random variable whose probability density function $f : \mathbb{R} \rightarrow [0, 1]$ is given by

$$f(x) = \begin{cases} \frac{2}{\pi R^2} \sqrt{R^2 - x^2} & \text{if } -R < x < R \\ 0 & \text{otherwise} \end{cases}$$

Find the distribution function of X . The distribution of X is called the semicircular distribution (see Figure 5.6).

Ex. 5.2.7. Let X be a random variable whose distribution function $F : \mathbb{R} \rightarrow [0, 1]$ is given by

$$F(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ \frac{2}{\pi} \arcsin(\sqrt{x}) & \text{if } 0 < x < 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

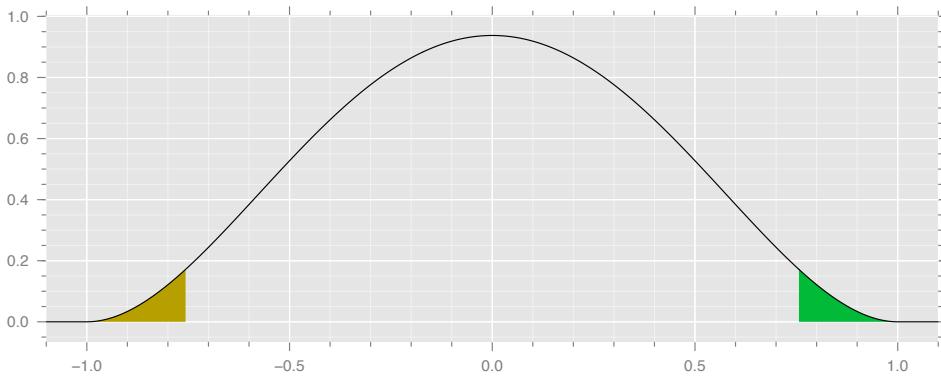


Figure 5.7: Computation of probabilities as area under the density curve. For symmetric distributions, it is enough to know the (cumulative) distribution function for positive values.

Find the probability density function of X . The distribution of X is called the standard arcsine law.

Ex. 5.2.8. Let X be a continuous random variable with probability density function f and distribution function F . Suppose f is a symmetric function, i.e. $f(x) = f(-x)$ for all $x \in \mathbb{R}$. Then show that

- (a) $P(X \leq 0) = P(X \geq 0) = \frac{1}{2}$,
- (b) for $x \geq 0$, $F(x) = \frac{1}{2} + P(0 \leq X \leq x)$,
- (c) for $x \leq 0$, $F(x) = P(X \geq -x) = \frac{1}{2} + P(0 \leq X \leq -x)$.

We have observed this fact for the normal distribution earlier (see Figure 5.7).

Ex. 5.2.9. Let $X \sim \text{Exp}(\lambda)$. The “90th percentile” is a value a such that X is larger than a 90% of the time. Find the 90th percentile of X by determining the value of a for which $P(X < a) = 0.9$.

Ex. 5.2.10. The “median” of a continuous random variable X is the value of x for which $P(X > x) = P(X < x) = \frac{1}{2}$.

- (a) If $X \sim \text{Uniform}(a, b)$ calculate the median of X .
- (b) If $Y \sim \text{Exp}(\lambda)$ calculate the median of Y .
- (c) Let $Z \sim \text{Normal}(\mu, \sigma^2)$. Show that the median of Z is μ .

Ex. 5.2.11. Let $X \sim \text{Normal}(\mu, \sigma^2)$. Show that $P(|X - \mu| < k\sigma)$ does not depend on the values of μ or σ . (Hint: Use a change of variables for the appropriate integral).

Ex. 5.2.12. Above we saw that exponential random variables satisfied the memoryless property, (5.2.6). It can be shown that any positive, continuous random variable with the memoryless property must be exponential. Follow the steps below to prove a slightly weakened version of this result. For all parts, suppose X is a positive, continuous random variable with the memoryless property for which the distribution function $F_X(t)$ has a continuous derivative for $t > 0$. Suppose further that $\lim_{t \rightarrow 0^+} F'(t)$ exists and call this quantity α . Let $G(t) = 1 - F_X(t) = P(X > t)$ and do the following.

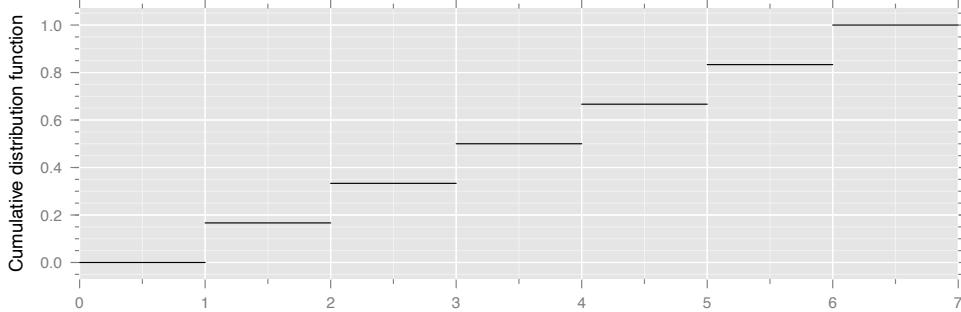


Figure 5.8: The cumulative distribution function for Exercise 5.2.14.

- (a) Use the memoryless property to show that $G(s+t) = G(s) \cdot G(t)$ for all positive s and t .
- (b) Use part (a) to conclude that $G'(t) = -\alpha G(t)$. (Hint: Take a derivative with respect to s and then take an appropriate limit).
- (c) It is a fact (which you may take as granted) that the differential equation from (b) has solutions of the form $G(t) = Ce^{-\alpha t}$. Use the fact that X is positive to explain why it must be that $C = 1$.
- (d) Use part (c) to calculate $F_X(t)$ and then differentiate to find $f_X(t)$.
- (e) Conclude that X must be exponentially distributed and determine the associated parameter in terms of α .

Ex. 5.2.13. Let X be a random variable with density $f(x) = 2x$ for $0 < x < 1$ (and $f(x) = 0$ otherwise). Calculate the distribution function of X .

Ex. 5.2.14. Let $X \sim \text{Uniform}(\{1, 2, 3, 4, 5, 6\})$. Despite the fact this is a discrete random variable without a density, the distribution function $F_X(x)$ is still defined. Find a piecewise defined expression for $F_X(x)$ (see Figure 5.8 for a plot).

Ex. 5.2.15. Suppose $F : \mathbb{R} \rightarrow [0, 1]$ is given by (5.2.3). Then show that

1. F is a monotonically increasing function.
2. $\lim_{x \rightarrow \infty} F(x) = 1$.
3. $\lim_{x \rightarrow -\infty} F(x) = 0$.
4. if, in addition, F is given by (5.2.4) then F is continuous.

Ex. 5.2.16. We use the notation as in Theorem 5.2.10.

- (a) Let

$$A_n = \left\{ k : 0 \leq k \leq n, np + a\sqrt{np(1-p)} \leq k \leq np + a\sqrt{np(1-p)} \right\}.$$

Show that

$$P(a \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq b) = \sum_{k \in A_n} P(S_n = k).$$

(b) Let

$$\zeta_{k,n} = \frac{k - np}{\sqrt{np(1-p)}}.$$

Using the definition of the Riemann integral show that

$$\lim_{n \rightarrow \infty} \sum_{k \in A_n} \frac{e^{-\frac{\zeta_{k,n}^2}{2}}}{\sqrt{2\pi np(1-p)}} = \int_a^b \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} dx$$

(c) Using Stirling's approximation show that

$$\lim_{n \rightarrow \infty} \sup_{k \in A_n} \frac{\binom{n}{k} p^k (1-p)^{n-k}}{\sqrt{2\pi np(1-p)} e^{-\frac{\zeta_{k,n}^2}{2}}} = 1$$

(d) Prove Theorem 5.2.10 by observing

$$\begin{aligned} P(a \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq b) &= \\ \sum_{k \in A_n} \frac{e^{-\frac{\zeta_{k,n}^2}{2}}}{\sqrt{2\pi np(1-p)}} + \sum_{k \in A_n} \frac{e^{-\frac{\zeta_{k,n}^2}{2}}}{\sqrt{2\pi np(1-p)}} &\left(\frac{\binom{n}{k} p^k (1-p)^{n-k}}{\frac{1}{\sqrt{2\pi np(1-p)}} e^{-\frac{\zeta_{k,n}^2}{2}}} - 1 \right) \end{aligned}$$

5.3 TRANSFORMATION OF CONTINUOUS RANDOM VARIABLES

In Section 3.3 we have discussed functions of discrete random variables and how to find their distributions. Suppose $g : \mathbb{R} \rightarrow \mathbb{R}$ and $Y = g(X)$, to find the distribution of Y we converted events associated with Y with events of X by inverting the function g . In the setting of continuous random variables distribution functions are used for calculating probabilities associated with functions of a known random variable. We next present a simple example for which $g(x) = x^2$ followed by a result that covers situations when $g(x)$ is any linear function.

EXAMPLE 5.3.1. Let $X \sim \text{Uniform}(0, 1)$ and let $Y = X^2$. What is the density for Y ?

Since X takes values on $(0, 1)$ and since $Y = X^2$, it will also be the case that Y takes values on $(0, 1)$. However, though X is uniform on the interval, there should be no expectation that Y will also be uniform. In fact, since squaring a positive number less than one results in a smaller number than the original, it should seem intuitive that results of Y will be more likely to be near to zero than they are to be near to one.

It is not easy to see how to calculate the density of Y directly from the density of X . However, it is a much easier task to compute the distribution of Y from the distribution of X . Therefore we will use the following plan in the calculation below – integrate $f_X(x)$ to find $F_X(x)$; use $F_X(x)$ to determine $F_Y(y)$; then differentiate $F_Y(y)$ to calculate $f_Y(y)$.

For the first step, note

$$F_X(x) = \int_{-\infty}^x f_X(x) dx = \begin{cases} 0 & \text{if } 0 < x \\ x & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$

Next, since Y takes values in $(0, 1)$, if $y \leq 0$ then $F_Y(y) = P(Y \leq y) = 0$. But if $y > 0$ then

$$F_Y(y) = P(Y \leq y) = P(X^2 \leq y) = P(-\sqrt{y} \leq X \leq \sqrt{y}).$$

Since X is always positive, the event $(X < -\sqrt{y})$ has zero probability we may connect this to the distribution of X by writing

$$\begin{aligned} F_Y(y) &= P(-\sqrt{y} \leq X \leq \sqrt{y}) \\ &= P(X < -\sqrt{y}) + P(-\sqrt{y} \leq X \leq \sqrt{y}) \\ &= P((X < -\sqrt{y}) \cup (-\sqrt{y} \leq X \leq \sqrt{y})) \\ &= P(X \leq \sqrt{y}) = F_X(\sqrt{y}). \end{aligned}$$

Therefore,

$$F_Y(y) = \begin{cases} 0 & \text{if } 0 \leq y \\ \sqrt{y} & \text{if } 0 < y < 1 \\ 1 & \text{if } y \geq 1 \end{cases}$$

and finally by using the fact that $F'(y) = f(y)$ we can determine that

$$f_Y(y) = \begin{cases} \frac{1}{2\sqrt{y}} & \text{if } 0 < y < 1 \\ 0 & \text{otherwise} \end{cases}$$

As noted in the beginning of this example, this distribution is far from uniform and gives much more weight to intervals close to zero than it does intervals close to one. ■

LEMMA 5.3.2. *Let $a \neq 0$ and $b \in \mathbb{R}$. Suppose X is a continuous random variable with probability density function f_X . Let $g(x) = ax + b$ be any non-constant linear function (so $a \neq 0$) and let $Y = g(X)$ then Y is also a continuous random variable whose density function f_Y is given by*

$$f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right), \quad (5.3.1)$$

for all $y \in \mathbb{R}$.

Proof- Let $y \in \mathbb{R}$. Assume first that $a > 0$. Then

$$P(Y \leq y) = P(aX + b \leq y) = P(X \leq \frac{y-b}{a}) = \int_{-\infty}^{\frac{y-b}{a}} f_X(z) dz$$

By a simple change of variable $z = \frac{u-b}{a}$ we obtain that

$$P(Y \leq y) = \int_{-\infty}^y \frac{1}{a} f_X\left(\frac{u-b}{a}\right) du. \quad (5.3.2)$$

If $a < 0$ then

$$P(Y \leq y) = P(aX + b \leq y) = P(X \geq \frac{y-b}{a}) = \int_{\frac{y-b}{a}}^{\infty} f_X(z) dz$$

Again a simple change of variable $z = \frac{u-b}{a}$, with $a < 0$, we obtain that

$$P(Y \leq y) = \int_{-\infty}^y \frac{1}{-a} f_X\left(\frac{u-b}{a}\right) du. \quad (5.3.3)$$

Using (5.3.2) and (5.3.3) we have that Y is a continuous random variable with density as in (5.3.1). ■

Lemma 5.3.2 provides a method to standardise the normal random variable.

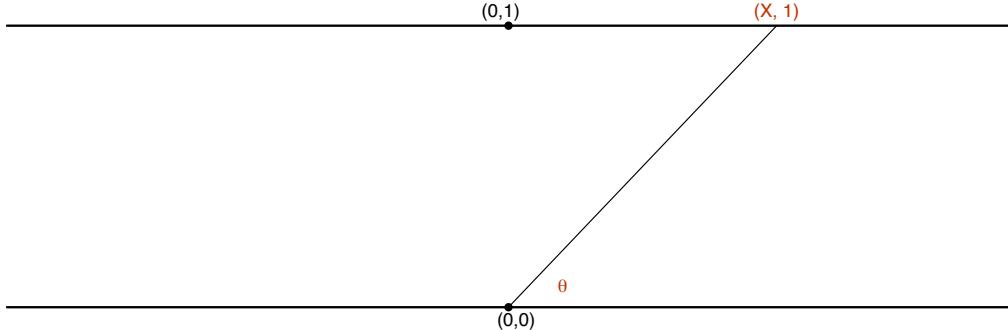


Figure 5.9: Illustration of Example 5.3.4.

COROLLARY 5.3.3. (a) Let $X \sim \text{Normal}(0, 1)$ and let $Y = aX + b$ with $a, b \in \mathbb{R}, a \neq 0$. Then, $Y \sim \text{Normal}(b, a^2)$.

(b) Let $X \sim \text{Normal}(\mu, \sigma^2)$ and let $Z = \frac{X-\mu}{\sigma}$. Then $Z \sim \text{Normal}(0, 1)$.

Proof - X has a probability density function given by (5.2.7).

(a) By Lemma 5.3.2, we have that the density of Y is given by

$$f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right) = \frac{1}{\sqrt{2\pi} |a|} e^{-\frac{(y-b)^2}{2a^2}},$$

for all $y \in \mathbb{R}$. Hence $Y \sim \text{Normal}(b, a^2)$.

(b) By Lemma 5.3.2, with $a = \frac{1}{\sigma}$ and $b = -\frac{\mu}{\sigma}$ we have that the density of Z is given by

$$f_Z(z) = \sigma f_X\left(\sigma(z + \frac{\mu}{\sigma})\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}},$$

for all $z \in \mathbb{R}$. Hence $Z \sim \text{Normal}(0, 1)$. ■

EXAMPLE 5.3.4. Consider the two parallel lines in \mathbb{R}^2 , given by $y = 0$ and $y = 1$. Piku is standing at the origin in the plane. She chooses an angle θ uniformly in $(0, \pi)$ and she draws a line segment between the lines $y = 0$ and $y = 1$ at an angle θ from the origin in \mathbb{R}^2 . Suppose the line segment meets the line $y = 1$ at the point $(X, 1)$. Find the probability density function of X .

First observe that $X = \tan(\frac{\pi}{2} - \theta)$. We shall first find the distribution function of X . Let $x \in \mathbb{R}$. Observe that $\tan(x)$ is a strictly increasing function in the interval $(-\frac{\pi}{2}, \frac{\pi}{2})$ and has an inverse denoted by $\arctan(x)$. So

$$\begin{aligned} P(X \leq x) &= P(\tan(\frac{\pi}{2} - \theta) \leq x) \\ &= P((\frac{\pi}{2} - \theta) \leq \arctan(x)) \\ &= P(\theta \geq \frac{\pi}{2} - \arctan(x)) \\ &= 1 - P(\theta \leq \frac{\pi}{2} - \arctan(x)) \end{aligned}$$

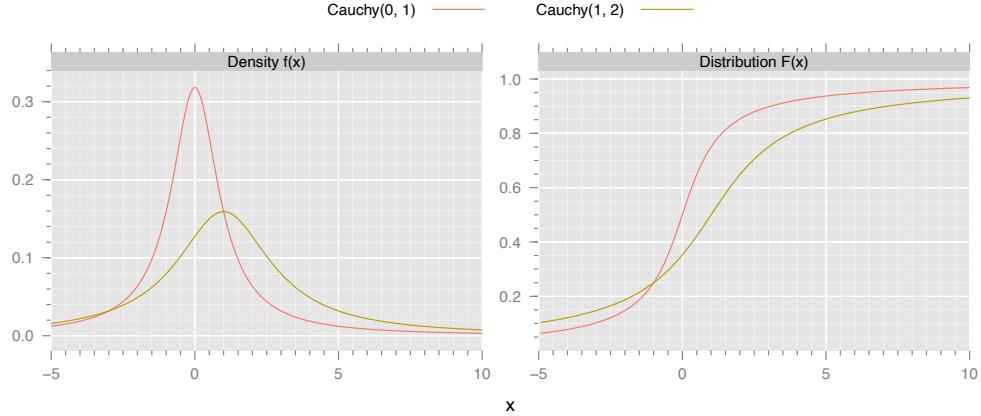


Figure 5.10: The shape of Cauchy density and cumulative distribution functions for selected parameter values.

For any $x \in \mathbb{R}$, $\frac{\pi}{2} - \arctan(x) \in (0, \pi)$. As θ has Uniform $(0, \pi)$ distribution, the above is

$$\begin{aligned} &= 1 - \frac{1}{\pi} \left(\frac{\pi}{2} - \arctan(x) \right) \\ &= \frac{1}{2} + \frac{1}{\pi} \arctan(x) \end{aligned}$$

Hence the distribution function of X is differentiable and therefore the probability density function of X is given by

$$f_X(x) = \frac{1}{\pi} \frac{1}{1+x^2},$$

for all $x \in \mathbb{R}$. Such a random variable is an example of a Cauchy distribution which we define more generally next. ■

DEFINITION 5.3.5. $X \sim \text{Cauchy}(\theta, \alpha^2)$: Let $\theta \in \mathbb{R}$ and let $\alpha > 0$. Then X is said to have a Cauchy distribution with parameters θ and α^2 if it has the density

$$f(x) = \frac{1}{\pi} \frac{\alpha}{\alpha^2 + (x - \theta)^2} \quad (5.3.4)$$

for all $x \in \mathbb{R}$. Here θ is referred to as the location parameter and α is referred to as the scale parameter. The distribution function of X is given by

$$F(x) = \frac{1}{\pi} \arctan\left(\frac{x - \theta}{\alpha}\right) \quad (5.3.5)$$

Figure 5.10 gives plots of the Cauchy density and distribution functions.

Similar computations as above are useful for simulations. Most computer programming languages and spreadsheets have a “Random” function designed to approximate a Uniform(0, 1) random variable. How could one use such a feature to simulate random variables with other densities? We start with an example.

EXAMPLE 5.3.6. If $X \sim \text{Uniform}(0, 1)$, our goal is to find a function $g : (0, 1) \rightarrow \mathbb{R}$ for which $Y = g(X) \sim \text{Exponential}(\lambda)$. We will try to find such a $g : (0, 1) \rightarrow \mathbb{R}$ which is strictly increasing so that it has an inverse. This will be important when it comes to relating the distributions of X and Y .

We require Y to $\text{Exponential}(\lambda)$. So the distribution function of Y is

$$F_Y(y) = \begin{cases} 0 & \text{if } y \leq 0 \\ 1 - e^{-\lambda y} & \text{if } y > 0 \end{cases}$$

But

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y) = P(X \leq g^{-1}(y))$$

where the final equality comes from our decree that the function g should be strictly increasing. Therefore,

$$F_Y(y) = F_X(g^{-1}(y)).$$

But the distribution function of a uniform random variable has previously been computed. Hence,

$$F_X(g^{-1}(y)) = \begin{cases} 0 & \text{if } g^{-1}(y) \leq 0 \\ g^{-1}(y) & \text{if } 0 < g^{-1}(y) < 1 \\ 1 & \text{if } g^{-1}(y) \geq 1 \end{cases}$$

Thus we are forced to have

$$g^{-1}(y) = 1 - e^{-\lambda y}$$

for $y > 0$. So inverting the above formula, we get $g : (0, 1) \rightarrow (0, \infty)$ is given by

$$g(x) = -\frac{1}{\lambda} \ln(1 - x),$$

for $x \in (0, 1)$. Hence,

$$X \sim \text{Uniform}(0, 1) \implies -\frac{1}{\lambda} \ln(1 - X) \sim \text{Exponential}(\lambda).$$

In conclusion one could view g as the inverse of F_Y , on $(0, \infty)$. It turns out that this is a general result. We state a special case of this in the lemma below. ■

LEMMA 5.3.7. Let $U \sim \text{Uniform}(0, 1)$ random variable. Let X be a continuous random variable such that its distribution function, F_X , is a strictly increasing continuous function. Then

(a) $Y = F_X^{-1}(U)$ has the same distribution as X .

(b) $Z = F_X(X)$ has the same distribution as U .

Proof- We observe that as F is strictly increasing continuous distribution function $F : \mathbb{R} \rightarrow (0, 1)$ and Range $(F) = (0, 1)$.

(a) We shall verify that Y and X have the same distribution function. Let $y \in \mathbb{R}$, then

$$F_Y(y) = P(Y \leq y) = P(F_X^{-1}(U) \leq y) = P(U \leq F_X(y)) = F_X(y)$$

Hence X and Y have the same distribution.

(b) We shall verify that Z and U have the same distribution function. Let $z \in \mathbb{R}$. If $z \leq 0$ then

$$P(Z \leq z) = P(F(X) \leq z) = 0$$

as $F : \mathbb{R} \rightarrow (0, 1)$. If $z \geq 1$ then

$$P(Z \leq z) = P(F(X) \leq z) = 1$$

as $F : \mathbb{R} \rightarrow (0, 1)$. If $0 < z < 1$ then $F^{-1}(z)$ is well defined as Range $(F) = (0, 1)$ and

$$P(Z \leq z) = P(F(X) \leq z) = P(X \leq F^{-1}(z)) = F(F^{-1}(z)) = z.$$

Hence Z and U have the same distribution. ■

The previous lemma may be generalized even to the case when F is not strictly increasing. It requires a concept called the generalized inverse. The interested reader will find it discussed in Exercise 5.3.12.

EXERCISES

Ex. 5.3.1. Let $X \sim \text{Uniform}(0, 1)$ and let $Y = \sqrt{X}$. Determine the density of Y .

Ex. 5.3.2. Let $X \sim \text{Uniform}(0, 1)$ and let $Z = \frac{1}{X}$. Determine the density of Z .

Ex. 5.3.3. Let $X \sim \text{Uniform}(0, 1)$. Let $r > 0$ and define $Y = rX$. Show that Y is uniformly distributed on $(0, r)$.

Ex. 5.3.4. Let $X \sim \text{Uniform}(0, 1)$. Let $Y = 1 - X$. Show that $Y \sim \text{Uniform}(0, 1)$ as well.

Ex. 5.3.5. Let $X \sim \text{Uniform}(0, 1)$. Let a and b be real numbers with $a < b$ and let $Y = (b - a)X + a$. Show that $Y \sim \text{Uniform}(a, b)$.

Ex. 5.3.6. Let $X \sim \text{Uniform}(0, 1)$. Find a function $g(x)$ (which is strictly increasing) such that the random variable $Y = g(X)$ has density $f_Y(y) = 3y^2$ for $0 < y < 1$ (and $f_Y(y) = 0$ otherwise).

Ex. 5.3.7. Let $X \sim \text{Normal}(\mu, \sigma^2)$. Let $g : (-\infty, \infty) \rightarrow \mathbb{R}$ be given by $g(x) = x^2$. Find the probability density function of $Y = g(X)$.

Ex. 5.3.8. Let $\alpha > 0$ and X be a random variable with the p.d.f given by

$$f(x) = \begin{cases} \frac{\alpha}{x^{\alpha+1}} & 1 \leq x < \infty \\ 0 & \text{otherwise} \end{cases}$$

The random variable X is said to have Pareto (α) distribution (see Figure 5.11).

- (a) Find the distribution of $X_1 = X^2$
- (b) Find the distribution of $X_2 = \frac{1}{X}$
- (c) Find the distribution of $X_3 = \ln(X)$

In the above exercises we assume that the transformation function is defined as above when the p.d.f of X is positive and zero otherwise.

Ex. 5.3.9. Let X be a continuous random variable with probability density function $f_X : \mathbb{R} \rightarrow \mathbb{R}$. Let $a > 0$, $b \in \mathbb{R}$ $Y = \frac{1}{a}(X - b)^2$. Show that Y is also a continuous random variable with probability density function $f_Y : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$f_Y(y) = \frac{\sqrt{a}}{2\sqrt{y}} [f_X(\sqrt{ay} + b) + f_X(-\sqrt{ay} + b)]$$

for $y > 0$.

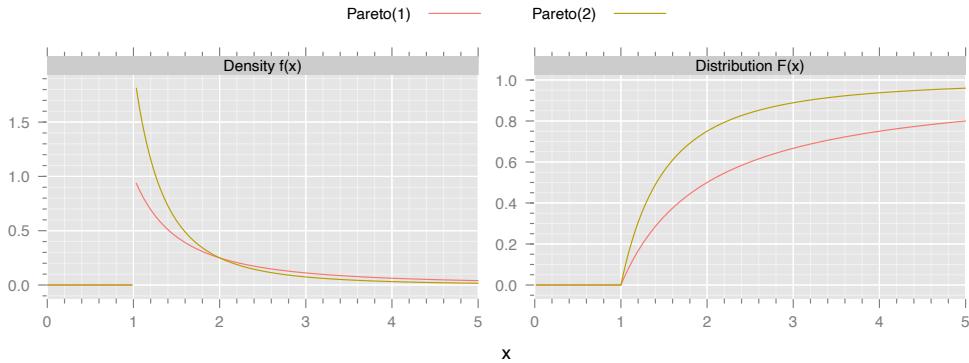


Figure 5.11: The shape of the pareto density and cumulative distribution functions.

Ex. 5.3.10. Let $-\infty \leq a < b \leq \infty$ and $I = (a, b)$ and $g : I \rightarrow \mathbb{R}$. Let X be a continuous random variable whose density f_X is zero on the complement of I . Set $Y = g(X)$.

(a) Let g be a differentiable strictly increasing function.

- (i) Show that inverse of g exists and g^{-1} is strictly increasing on $g(I)$.
- (ii) For any $y \in \mathbb{R}$, show that $P(Y \leq y) = P(X \leq g^{-1}(y))$
- (iii) Show that Y has a density $f_Y(\cdot)$ given by

$$f_Y(y) = f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y).$$

(b) Let g be a differentiable strictly decreasing function.

- (i) Show that inverse of g exists and g^{-1} is strictly decreasing on $g(I)$.
- (ii) For any $y \in \mathbb{R}$, show that $P(Y \leq y) = 1 - P(X \leq g^{-1}(y))$
- (iii) Show that Y has a density $f_Y(\cdot)$ given by

$$f_Y(y) = f_X(g^{-1}(y)) \left(-\frac{d}{dy} g^{-1}(y) \right).$$

Ex. 5.3.11. Let X be a random variable having an exponential density. Let $g : [0, \infty) \rightarrow \mathbb{R}$ be given by $g(x) = x^{\frac{1}{\beta}}$, for some $\beta \neq 0$. Find the probability density function of $Y = g(X)$.

Ex. 5.3.12. Let $U \sim \text{Uniform}(0, 1)$. Let X be a continuous random variable with a distribution function F . Extend $F : \mathbb{R} \rightarrow \mathbb{R}$ to $F : \mathbb{R} \cup \{-\infty\} \cup \{\infty\} \rightarrow \mathbb{R}$ by setting $F(\infty) = 1$ and $F(-\infty) = 0$. Define the generalised inverse of F , $G : [0, 1] \rightarrow \mathbb{R} \cup \{-\infty\} \cup \{\infty\}$ by

$$G(y) = \inf\{x \in \mathbb{R} : F(x) \geq y\}.$$

Show that

- (a) Show that for all $y \in [0, 1]$, $F(G(y)) = y$.
- (b) Show that for all $x \in \mathbb{R}$ and $y \in [0, 1]$

$$F(x) \geq y \iff x \geq G(y).$$

- (c) $Y = G(U)$ has the same distribution as X .
- (d) $Z = F(X)$ has the same distribution as U .

5.4 MULTIPLE CONTINUOUS RANDOM VARIABLES

When analyzing multiple random variables at once, one may consider a “joint density” analogous to the joint distribution of the discrete variable case. In this section we will restrict considerations to only two random variables, but we shall see in Chapter 8 that the definitions and results all generalize to any finite collection of variables.

THEOREM 5.4.1. *Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a non-negative function, piecewise-continuous in each variable for which*

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1.$$

For a Borel set $A \subset \mathbb{R}^2$ define

$$P(A) = \int_A f(x, y) dx dy.$$

Then P is a probability on \mathbb{R}^2 and f is called the density for P .

Proof- The proof of the theorem is essentially the same as in the one-variable version of Theorem 5.1.5. We will not reproduce it here. As in the discrete case we will typically associate such densities with random variables. ■

DEFINITION 5.4.2. *A pair of random variables (X, Y) is said to have a joint density $f(x, y)$ if for every Borel set $A \subset \mathbb{R}^2$*

$$P((X, Y) \in A) = \int_A f(x, y) dx dy.$$

As in the one-variable case we describe this in terms of “Borel sets” to be precise, but in practice we will only consider sets A which are simple regions in the plane. In fact regions such as $(-\infty, a] \times (-\infty, b]$, for all real numbers a, b are enough to characterise the joint distribution. As in the one variable case we can define a “joint distribution function” of (X, Y) as

$$F_{(X,Y)}(a, b) = P((X \leq a) \cap (Y \leq b)) = \int_{-\infty}^a \int_{-\infty}^b f(z, w) dw dz \quad (5.4.1)$$

for all $a, b \in \mathbb{R}$. We will usually denote the joint distribution function by F omitting the subscripts unless it is particularly needed. One can state and prove a similar type of result as Theorem 5.2.5 for $F(a, b)$ when (X, Y) have a joint density. In particular, we can conclude that since the joint densities are assumed to be piecewise continuous, the corresponding distribution functions are piecewise differentiable. Further, the joint distribution of two continuous random variables (X, Y) are completely determined by their joint distribution function F . That is, if we know the value of $F(a, b)$ for all $a, b \in \mathbb{R}$, we could use multivariable calculus to differentiate $F(a, b)$ to find $f(a, b)$. Then $P((X, Y) \in A)$ for any event A is obtained by integrating the joint density f over the event A . We illustrate this with a couple of examples.

EXAMPLE 5.4.3. Consider the open rectangle in \mathbb{R}^2 given by $R = (0, 1) \times (3, 5)$ and $|R| = 2$ denote its area. Let (X, Y) have a joint density $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by

$$f(x, y) = \begin{cases} \frac{1}{2} & \text{if } (x, y) \in R \\ 0 & \text{otherwise.} \end{cases}$$

The above is clearly a density function. So for any rectangle $A = (a, b) \times (c, d) \subset R$,

$$P((X, Y) \in A) = \int_a^b \int_c^d f(x, y) dx dy = \frac{(b-a)(d-c)}{2} = \frac{|A|}{|R|}.$$

■

In general one can use the following definition to define a uniform random variable on the plane.

DEFINITION 5.4.4. Let $D \subset \mathbb{R}^2$ be non-empty and with positive area (assume D is a Borel set or in particular f or any simple region whose area is well defined). Then $(X, Y) \sim \text{Uniform}(D)$ if it has a joint probability density function given by $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by

$$f(x, y) = \begin{cases} \frac{1}{|D|} & \text{if } (x, y) \in D \\ 0 & \text{otherwise,} \end{cases}$$

where $|D|$ denotes the area of D .

When $(X, Y) \sim \text{Uniform}(D)$ then the probability that (X, Y) lies in a region $A \subset D$ is proportional to the area of A .

EXAMPLE 5.4.5. Let (X, Y) have a joint density $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by

$$f(x, y) = \begin{cases} x + y & \text{if } 0 < x < 1, 0 < y < 1 \\ 0 & \text{otherwise} \end{cases}$$

We note that this really does describe a density. The function $f(x, y)$ is non-negative and

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy &= \int_0^1 \int_0^1 x + y dx dy \\ &= \int_0^1 \left(\frac{1}{2}x^2 + xy \right) \Big|_{x=0}^{x=1} dy \\ &= \int_0^1 \frac{1}{2} + y dy \\ &= \frac{1}{2}y + \frac{1}{2}y^2 \Big|_{y=0}^{y=1} = 1. \end{aligned}$$

Calculating a probability such as $P((X < \frac{1}{2}) \cap (Y < \frac{1}{2}))$ requires integrating over the appropriate region.

$$\begin{aligned} P((X < \frac{1}{2}) \cap (Y < \frac{1}{2})) &= \int_{-\infty}^{1/2} \int_{-\infty}^{1/2} f(x, y) dx dy \\ &= \int_0^{1/2} \int_0^{1/2} x + y dx dy \\ &= \int_0^{1/2} \frac{1}{8} + \frac{1}{2}y dy \\ &= \frac{1}{8}. \end{aligned}$$

A probability only involving one variable may still be calculated from the joint density. For instance $P(X < \frac{1}{2})$ does not appear to involve Y , but this simply means that Y is unrestricted and the corresponding integral should range over all possible values of Y . Therefore,

$$\begin{aligned} P(X < \frac{1}{2}) &= \int_{-\infty}^{\infty} \int_{-\infty}^{1/2} f(x, y) dx dy \\ &= \int_0^1 \int_0^{1/2} x + y dx dy = \frac{3}{8}. \end{aligned}$$

It is just as easy to compute that $P(Y < \frac{1}{2}) = \frac{3}{8}$. Note that these computations also demonstrate that X and Y are not independent since

$$P(X < \frac{1}{2}) \cdot P(Y < \frac{1}{2}) = \frac{9}{64} \neq P((X < \frac{1}{2}) \cap (Y < \frac{1}{2})).$$

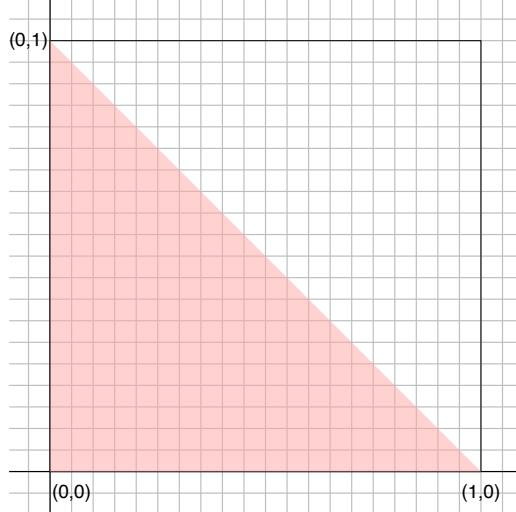


Figure 5.12: The subset A of the unit square that represents the region $x + y < 1$.

A probability such as $P(X + Y < 1)$ can be found by integrating over a non-rectangular region in the plane, as shown in Figure 5.12. Let $A = \{(x, y) | x + y < 1\}$. Then

$$\begin{aligned} P(X + Y < 1) &= \int_A f(x, y) dx dy \\ &= \int_0^1 \int_0^{1-y} x + y dx dy \\ &= \int_0^1 \frac{1}{2}x^2 + xy \Big|_0^{1-y} dy \\ &= \int_0^1 \frac{1}{2}(1-y)^2 + (1-y)y dy \\ &= \int_0^1 \frac{1}{2} - \frac{1}{2}y^2 dy \\ &= \frac{1}{3}. \end{aligned}$$

■

5.4.1 Marginal Distributions

As in the discrete case, when we begin with the joint density of many random variables, but want to speak of the distribution of an individual variable we will frequently refer to it as a “marginal distribution”.

Suppose (X, Y) are random variables and have a joint probability density function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Then we observe that

$$P(X \leq x) = P(X \leq x, -\infty < Y < \infty) = \int_{-\infty}^x \int_{-\infty}^{\infty} f(u, y) dy du.$$

If $g : \mathbb{R} \rightarrow \mathbb{R}$ is given by

$$g(u) = \int_{-\infty}^{\infty} f(u, y) dy$$

then

$$P(X \leq x) \int_{-\infty}^x g(u) du.$$

Using Theorem 5.2.5, by the continuity assumptions on f , we find that the random variable X is also a continuous random variable with probability density function of X given by

$$f_X(x) = g(x) = \int_{-\infty}^{\infty} f(x, y) dy. \quad (5.4.2)$$

As it was derived from a joint probability density function, the density of X is referred to as the marginal density of X . Similarly one can show that Y is also a continuous random variable and its marginal density is given by

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx. \quad (5.4.3)$$

EXAMPLE 5.4.6. (Example 5.4.3 contd.) Going back to Example 5.4.3, we can compute the marginal density of X and Y . The marginal density of X is given by

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy = \begin{cases} \int_3^5 \frac{1}{2} & \text{if } 0 < x < 1 \\ 0 & \text{otherwise.} \end{cases} = \begin{cases} 1 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise.} \end{cases}$$

The marginal density of Y is given by

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx = \begin{cases} \int_0^1 \frac{1}{2} & \text{if } 3 < y < 5 \\ 0 & \text{otherwise.} \end{cases} = \begin{cases} \frac{1}{2} & \text{if } 3 < y < 5 \\ 0 & \text{otherwise.} \end{cases}$$

So we observe that $X \sim \text{Uniform}(0, 1)$ and $Y \sim \text{Uniform}(3, 5)$. ■

While it is routine to find the marginal densities from the joint density there is no standard way to get to the joint from the marginals. Part of the reason for this difficulty is that the marginal densities offer no information about how the variables relate to each other, which is critical information for determining how they behave jointly. However, in the case that the random variables happen to be independent there is a convenient relationship between the joint and marginal densities.

5.4.2 Independence

THEOREM 5.4.7. *Let f be the joint density of random variables X and Y and let f_X and f_Y be the respective marginal densities. Then*

$$f(x, y) = f_X(x)f_Y(y)$$

if and only if X and Y are independent.

Proof - First suppose X and Y are independent and consider the quantity $P((X \leq x) \cap (Y \leq y))$. On one hand independence gives

$$P((X \leq x) \cap (Y \leq y)) = P(X \leq x)P(Y \leq y) = F_X(x)F_Y(y) \quad (5.4.4)$$

On the other hand, integrating the joint density yields

$$P((X \leq x) \cap (Y \leq y)) = \int_{-\infty}^x \int_{-\infty}^y f(x, y) dx dy. \quad (5.4.5)$$

Since equations 5.4.4 and 5.4.5 are equal we may differentiate both with respect to each of the variables x and y and they remain equal. However, differentiating the former gives $f_X(x)f_Y(y)$ because of the relationship between the distribution and the density, while differentiating the latter yields $f(x, y)$ by a two-fold application of the fundamental theorem of calculus.

To prove the opposite direction, suppose $f(x, y) = f_X(x)f_Y(y)$. Let A and B be Borel sets in \mathbb{R} . Then

$$\begin{aligned} P((X \in A) \cap (Y \in B)) &= \int_B \int_A f(x, y) dx dy \\ &= \int_B \int_A f_X(x)f_Y(y) dx dy \\ &= \left(\int_A f_X(x) dx \right) \left(\int_B f_Y(y) dy \right) \\ &= P(X \in A)P(Y \in B) \end{aligned}$$

Since this is true for all sets such sets A and B , the variables X and Y are independent. ■

EXAMPLE 5.4.8. (Example 5.4.3 contd.) We had observed that if $(X, Y) \sim \text{Uniform}(R)$ then $X \sim \text{Uniform}(0, 1)$ and $Y \sim \text{Uniform}(3, 5)$. Note further that

$$f(x, y) = f_X(x)f_Y(y)$$

for all $x, y \in R$. Consequently X, Y are independent as well. ■

It is tempting to generalise and say that $(X, Y) \sim \text{Uniform}(D)$ for a region D with non-trivial area then X and Y would be independent. This is not the case, we illustrate in the example below.

EXAMPLE 5.4.9. Consider the open disk in \mathbb{R}^2 given by $C = \{(x, y) : x^2 + y^2 < 25\}$ and $|C| = 25\pi$ denote its area. Let (X, Y) have a joint density $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by

$$f(x, y) = \begin{cases} \frac{1}{|C|} & \text{if } (x, y) \in C \\ 0 & \text{otherwise.} \end{cases}$$

As before for any Borel $A \subset C$,

$$P((X, Y) \in A) = \frac{|A|}{|C|},$$

and the probability that (X, Y) lies in A is proportional to the area of A . However the marginal density calculation is a little different. The marginal density of X is given by

$$\begin{aligned} f_X(x) &= \int_{-\infty}^{\infty} f(x, y) dy = \begin{cases} \int_{-\sqrt{25-x^2}}^{\sqrt{25-x^2}} \frac{1}{|C|} dy & \text{if } -5 < x < 5 \\ 0 & \text{otherwise.} \end{cases} \\ &= \begin{cases} \frac{2}{25\pi} \sqrt{25-x^2} & \text{if } -5 < x < 5 \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

The distribution of X is the Semi-circular law described in Exercise 5.2.6. As the joint density f is symmetric in x and y (i.e $f(x, y) = f(y, x)$) the marginal density of Y is the same as that of X (why?). It is easy to see

$$\frac{1}{25\pi} = f(0, 0) \neq f_X(0)f_Y(0) = \frac{4}{25\pi^2}$$

Consequently X, Y are not independent. This fact should make intuitive sense as well, for if X happens to take a value near 5 or -5 the range of possible values of Y is much more restricted than if X takes a value near 0. ■

We shall see the utility of independence when computing distributions of various functions of independent random variables (see Section 5.5). Independence of random variables also makes it easier to compute their joint density and hence probabilities. For instance, consider the following example.

EXAMPLE 5.4.10. Suppose $X \sim \text{Exponential}(\lambda_1)$, $Y \sim \text{Exponential}(\lambda_2)$ are independent random variables. Find $P(X - Y < 0)$.

The joint density of (X, Y) is given by

$$f(x, y) = f_X(x)f_Y(y) = \begin{cases} \lambda_1 \lambda_2 e^{-(\lambda_1 x + \lambda_2 y)} & \text{if } x > 0 \text{ and } y > 0 \\ 0 & \text{otherwise} \end{cases}$$

Therefore

$$\begin{aligned}
 P(X - Y < 0) &= \int_0^\infty \int_0^y \lambda_1 \lambda_2 e^{-(\lambda_1 x + \lambda_2 y)} dx dy = \lambda_1 \lambda_2 \int_0^\infty e^{-\lambda_2 y} \left[\int_0^y e^{-\lambda_1 x} dx \right] dy \\
 &= \lambda_1 \lambda_2 \int_0^\infty e^{-\lambda_2 y} \frac{1}{\lambda_1} [1 - e^{-\lambda_1 y}] dy \\
 &= \lambda_2 \left[\int_0^\infty e^{-\lambda_2 y} - e^{-(\lambda_1 + \lambda_2)y} dy \right] \\
 &= \lambda_2 \left[\frac{-1}{\lambda_2} (e^{-\lambda_2 y} |_0^\infty) + \frac{1}{\lambda_1 + \lambda_2} (e^{-(\lambda_1 + \lambda_2)y} |_0^\infty) \right] \\
 &= \lambda_2 \left[\frac{1}{\lambda_2} - \frac{1}{\lambda_1 + \lambda_2} \right] \\
 &= \frac{\lambda_1}{\lambda_1 + \lambda_2}.
 \end{aligned}$$

Similarly one can also compute $P(Y - X < 0) = \frac{\lambda_2}{\lambda_1 + \lambda_2}$. This fact is quite useful when using exponential random variables to model waiting times, for $P(X - Y < 0) = P(X < Y)$, so we have determined the probability that one waiting time will be shorter than another. ■

5.4.3 Conditional Density

In Section 3.2.2 we have seen the notion of conditional distributions for discrete random variables and in Section 4.4 we have seen the notions of conditional expectation and variance for discrete random variables. Suppose X measures the parts per million of a particulate matter less than 10 microns in the air and Y is the incidence rate of asthma in the population. It is clear that X and Y ought to be related; for the distribution of one affects the other. Towards this, in this section we shall discuss conditional distributions for two continuous random variables having a joint probability density function. We recall from Definition 3.2.5 that if X is a random variable on a sample space S and $A \subset S$ be an event such that $P(A) > 0$, then the probability Q described by

$$Q(B) = P(X \in B | A)$$

is called the conditional distribution of X given the event A .

Suppose X and Y have a joint probability density function f . Given our discussion for discrete random variables it is natural to characterise the conditional distribution of X given some information on Y . In the discrete setting we typically considered an event $A = \{Y = b\}$ for some real number b in the range of Y . In the continuous setting such an event A would have zero probability, so the usual way of conditioning on an event would not be possible. However, there is a way to make such a conditioning meaningful and precise provided $f_Y(b) > 0$, where f_Y is the marginal density of Y .

Suppose we wish to find the following :

$$P(X \in [3, 4] | Y = b).$$

We shall argue heuristically and arrive at an expression for the above probability. Suppose the marginal density of X is $f_X(\cdot)$, and that of Y is $f_Y(\cdot)$. Assume first that f_Y is piecewise continuous and $f_Y(b) > 0$. Then it is a standard fact from real analysis to see that

$$P(Y \in [b, b + \frac{1}{n})) > 0,$$

for all $n \geq 1$. One can then view the conditional probability as before, that is

$$\begin{aligned} P(X \in [3, 4] \mid X \in [b, b + \frac{1}{n})) &= \frac{P(X \in [3, 4] \cap X \in [b, b + \frac{1}{n}))}{P(X \in [b, b + \frac{1}{n}))} \\ &= \frac{\int_b^{b+\frac{1}{n}} f_3(u, v) du}{\int_b^{b+\frac{1}{n}} f_X(u) du} \\ &= \frac{\int_b^{b+\frac{1}{n}} \left(n \int_b^{b+\frac{1}{n}} f(u, v) du \right) dv}{n \int_b^{b+\frac{1}{n}} f_X(u) du} \end{aligned}$$

From facts in real analysis (under some mild assumptions on f) the following can be established,

$$\lim_{n \rightarrow \infty} n \int_b^{b+\frac{1}{n}} f(u, v) du = f(b, v),$$

for all real numbers v and

$$\lim_{n \rightarrow \infty} n \int_b^{b+\frac{1}{n}} f_X(u) du = f_X(b).$$

We have seen earlier (see Exercise 1.1.13 (b))

$$\lim_{n \rightarrow \infty} P(Y \in [b, b + \frac{1}{n})) = P(Y = b).$$

Hence it would be reasonable to argue that $P(X \in [3, 4] \mid Y = b)$ ought to be defined as

$$P(X \in [3, 4] \mid Y = b) = \frac{\int_b^4 f(b, v) dv}{f_Y(b)}.$$

With the above motivation we are now ready to define conditional densities for two random variables.

DEFINITION 5.4.11. Let (X, Y) be random variables having joint density f . Let the marginal density of Y be $f_Y(\cdot)$. Suppose b is a real number such that $f_Y(b) > 0$ and is continuous at b then conditional density of X given $Y = b$ is given by

$$f_{X|Y=b}(x) = \frac{f(x, b)}{f_Y(b)} \quad (5.4.6)$$

for all real numbers x . Similarly, let the marginal density of X be $f_X(\cdot)$. Suppose a is a real number such that $f_X(a) > 0$ and is continuous at a then conditional density of Y given $X = a$ is given by

$$f_{Y|X=a}(y) = \frac{f(a, y)}{f_X(a)}$$

for all real numbers y .

This definition genuinely defines a probability density function, for $f_{X|Y=b}(x) \geq 0$ since it is the ratio of a non-negative quantity and a positive quantity. Moreover,

$$\begin{aligned}\int_{-\infty}^{\infty} f_{X|Y=b}(x)dx &= \int_{-\infty}^{\infty} \frac{f(x,b)}{f_Y(b)} dx \\ &= \frac{1}{f_Y(b)} \int_{-\infty}^{\infty} f(x,b)dx = \frac{1}{f_Y(b)} f_Y(b) = 1\end{aligned}$$

Note that if X and Y are independent then

$$f_{X|Y=b}(x) = \frac{f(x,b)}{f_Y(b)} = \frac{f_X(x)f_Y(b)}{f_Y(b)} = f_X(x).$$

One can use the conditional density to compute the conditional probabilities, namely if (X, Y) are random variables having joint density f and b is a real number such that its marginal density has the property $f_Y(b) > 0$ then

$$P(X \in A | Y = b) = \int_A f_{X|Y=b}(x)dx = \int_A \frac{f(x,b)}{f_Y(b)} dx.$$

We conclude this section with two examples where we compute conditional densities. In both the examples the dependencies between the random variables imply that the conditional distributions are different from the marginal distributions.

EXAMPLE 5.4.12. Let (X, Y) have joint probability density function f given by

$$f(x,y) = \frac{\sqrt{3}}{4\pi} e^{-\frac{1}{2}(x^2 - xy + y^2)} \quad -\infty < x, y < \infty.$$

Let $x \in \mathbb{R}$, then the marginal density of X at x is given by

$$f_X(x) = \int_{-\infty}^{\infty} f(x,y)dy = \int_{-\infty}^{\infty} \frac{\sqrt{3}}{4\pi} e^{-\frac{1}{2}(x^2 - xy + y^2)} dy$$

By a standard completing the square computation, $\frac{1}{2}(x^2 - xy + y^2) = \frac{3x^2}{8} + \frac{1}{2}(y - \frac{x}{2})^2$. Therefore,

$$f_X(x) = \frac{\sqrt{3}}{4\pi} e^{-\frac{3x^2}{8}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(y - \frac{x}{2})^2} dy$$

Observing that $\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y - \frac{x}{2})^2} dy = 1$ (why ?), we have

$$f_X(x) = \frac{\sqrt{3}}{4\pi} e^{-\frac{3x^2}{8}} \sqrt{2\pi} = \sqrt{\frac{3}{4}} \frac{1}{\sqrt{2\pi}} e^{-\frac{3x^2}{8}}$$

Hence X is a Normal random variable with mean 0 and variance $\frac{4}{3}$. By symmetry (or calculating similarly as above) we can also show that Y is a Normal random variable with mean 0 and variance $\frac{4}{3}$. Also, we can easily see that

$$f_X(x)f_Y(y) = \frac{3}{8\pi} e^{-\frac{3}{8}(x^2 + y^2)} \neq \frac{\sqrt{3}}{4\pi} e^{-\frac{1}{2}(x^2 - xy + y^2)} = f(x,y)$$

for many $x, y \in \mathbb{R}$. Hence X and Y are not independent. Note that $f_X(x) \neq 0$ for all real numbers x and is continuous at all $x \in \mathbb{R}$. Fix $x \in \mathbb{R}$, the conditional density of Y given $X = x$ is given by

$$f_{Y|X=x}(y) = \frac{f(x,y)}{f_X(x)} = \frac{\frac{\sqrt{3}}{4\pi} e^{-\frac{1}{2}(x^2 - xy + y^2)}}{\sqrt{\frac{3}{4}} \frac{1}{\sqrt{2\pi}} e^{-\frac{3x^2}{8}}} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y - \frac{x}{2})^2} \quad \forall y \in \mathbb{R}.$$

Hence though the marginal distribution of Y is $\text{Normal}(0, \frac{4}{3})$, the conditional distribution of Y given $X = x$ is Normal with mean $\frac{x}{2}$ and variance 1. Put another way, if we are given that $X = x$ the mean of Y changes from 0 to x and the variance reduces from $\frac{4}{3}$ to 1.

Such a pair (X, Y) is an example of a bivariate normal random variable and will be discussed in detail in Section 6.4. ■

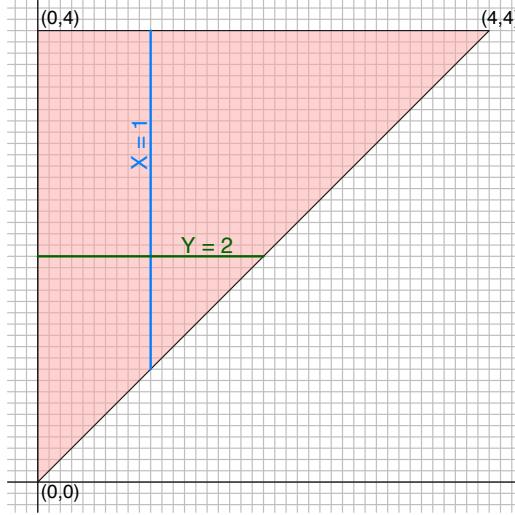


Figure 5.13: The region $T = \{(x, y) \mid 0 < x < y < 4\}$ from Example 5.4.13.

EXAMPLE 5.4.13. Suppose $T = \{(x, y) \mid 0 < x < y < 4\}$ and let $(X, Y) \sim \text{Uniform}(T)$. Therefore its joint density is given by (see Figure 5.13)

$$f(x, y) = \begin{cases} \frac{1}{8} & \text{if } (x, y) \in T \\ 0 & \text{otherwise.} \end{cases}$$

The marginal density of X is given by

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy = \begin{cases} \int_x^4 \frac{1}{8} dy & \text{if } 0 < x < 4 \\ 0 & \text{otherwise.} \end{cases} = \begin{cases} \frac{4-x}{8} & \text{if } 0 < x < 4 \\ 0 & \text{otherwise.} \end{cases}$$

The marginal density of Y is given by

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx = \begin{cases} \int_0^y \frac{1}{8} dx & \text{if } 0 < y < 4 \\ 0 & \text{otherwise.} \end{cases} = \begin{cases} \frac{y}{8} & \text{if } 0 < y < 4 \\ 0 & \text{otherwise.} \end{cases}$$

Let us fix $0 < b < 4$. So $f_Y(\cdot)$ is non-zero at b and is continuous at b . The conditional density of $(X \mid Y = b)$ is given by

$$f_{X|Y=b}(x) = \frac{f(x, b)}{f_Y(b)} = \begin{cases} \frac{1/8}{b/8} & \text{if } 0 < x < b \\ 0 & \text{otherwise.} \end{cases} = \begin{cases} \frac{1}{b} & \text{if } 0 < x < b \\ 0 & \text{otherwise.} \end{cases}$$

Therefore $(X \mid Y = b) \sim \text{Uniform}(0, b)$. Similarly if we fix $0 < a < 4$, we observe $f_X(\cdot)$ is non-zero at a and is continuous at a . The conditional density of $(Y \mid X = a)$ is given by

$$f_{Y|X=a}(y) = \frac{f(a, y)}{f_X(a)} = \begin{cases} \frac{1/8}{(4-a)/8} & \text{if } a < y < 4 \\ 0 & \text{otherwise.} \end{cases} = \begin{cases} \frac{1}{4-a} & \text{if } a < y < 4 \\ 0 & \text{otherwise.} \end{cases}$$

Therefore $(Y | X = a) \sim \text{Uniform}(a, 4)$.

Clearly X and Y are continuous random variables with distributions that are not uniform, but the conditional distributions turn out to be uniform. ■

EXERCISES

Ex. 5.4.1. Let (X, Y) be random variables whose probability density function is given by $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Find the probability density function of X and probability density function of Y in each of the following cases:-

- (a) $f(x, y) = (x + y)$ if $0 \leq x \leq 1, 0 \leq y \leq 1$ and 0 otherwise
- (b) $f(x, y) = 2(x + y)$ if $0 \leq x \leq y \leq 1$ and 0 otherwise
- (c) $f(x, y) = 6x^2y$ if $0 \leq x \leq 1, 0 \leq y \leq 1$ and 0 otherwise
- (d) $f(x, y) = 15x^2y$ if $0 \leq x \leq y \leq 1$ and 0 otherwise

Ex. 5.4.2. Let $c > 0$. Suppose that X and Y are random variables with joint probability density

$$f(x, y) = \begin{cases} c(xy + 1) & \text{if } 0 \leq x \leq 1 \text{ and } 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find c .
- (b) Compute the marginal densities $f_X(\cdot)$ and $f_Y(\cdot)$ and the conditional density $f_{X|Y=b}(\cdot)$

Ex. 5.4.3. Let $A = \{(x, y) \in \mathbb{R}^2 : x > 0, y > 0, x + y < 1\}$ and let X and Y be random variables defined by the joint density $f(x, y) = 24xy$ if $(x, y) \in A$ (and $f(x, y) = 0$ otherwise).

- (a) Verify the claim that $f(x, y)$ is a density.
- (b) Show that X and Y are dependent random variables.
- (c) Explain why (b) doesn't violate Theorem 5.4.7 despite the fact that $24xy$ is a product of a function of x with a function of y .

Ex. 5.4.4. Consider the set $D = [-1, 1] \times [-1, 1]$. Let

$$L = \{(x, y) \in D : x = 0 \text{ or } x = -1 \text{ or } x = 1 \text{ or } y = 0 \text{ or } y = 1 \text{ or } y = -1\}$$

be the lines that create a tiling of D . Suppose we drop a coin of radius R at a uniformly chosen point in D what is the probability that it will intersect the set L ?

Ex. 5.4.5. Let X and Y be two independent uniform $(0, 1)$ random variables. Let $U = \max(X, Y)$ and $V = \min(X, Y)$.

- (a) Find the joint distribution of U, V .
- (b) Find the conditional distribution of $(V | U = 0.5)$

Ex. 5.4.6. Suppose X is a random variable with density

$$f(x) = \begin{cases} cx^2(1-x) & \text{for } 0 \leq x \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Find:

- (a) the value of c .
- (b) the distribution function of X .
- (c) the conditional probability $P(X > 0.2 | X < 0.5)$.

Ex. 5.4.7. Suppose $g : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous probability density function, such that $g(x) = 0$ when $x \notin [0, 1]$. Let $D \subset \mathbb{R}^2$ be given by

$$D = \{(x, y) : x \in \mathbb{R} \text{ and } 0 \leq y \leq g(x)\}$$

Let (X, Y) be uniformly distributed on D . Find the probability density function of X .

Ex. 5.4.8. Continuous random variables X and Y have a joint density

$$f(x, y) = \begin{cases} \frac{1}{24}, & \text{for } 0 < x < 6, 0 < y < 4 \\ 0, & \text{elsewhere.} \end{cases}$$

- (a) Find $P(2Y > X)$.
- (b) Are X and Y independent?

Ex. 5.4.9. Let

$$f(x, y) = \begin{cases} \eta(y - x)^\gamma & \text{if } 0 \leq x < y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- (a) For what values of γ can η be chosen so that f be a joint probability density function of X , Y .
- (b) Given a γ from part (a), what is the value of η ?
- (c) Given a γ and η from parts (a) and (b), find the marginal densities of X and Y .

Ex. 5.4.10. Let $D = \{(x, y) : x^3 \leq y \leq x\}$. A point (X, Y) is chosen uniformly from D . Find the joint probability density function of X and Y .

Ex. 5.4.11. Let X and Y be two random variables with the joint p.d.f given by

$$f(x, y) = \begin{cases} ae^{-by} & 0 \leq x \leq y \\ 0 & \text{otherwise} \end{cases}$$

Find a conditions on a and b that make this a joint probability density function.

Ex. 5.4.12. Suppandi and Meera plan to meet at Gopalan Arcade between 7pm and 8pm. Each will arrive at a time (independent of each other) uniformly between 7pm and 8pm and will wait for 15 minutes for the other person before leaving. Find the probability that they will meet ?

5.5 FUNCTIONS OF INDEPENDENT RANDOM VARIABLES

In Section 5.3 we have seen how to compute the distribution of $Y = g(X)$ from the distribution of X for various $g : \mathbb{R}^2 \rightarrow \mathbb{R}$. Suppose (X, Y) are random variables having a joint probability density function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Let $h : \mathbb{R}^2 \rightarrow \mathbb{R}$. A natural follow up objective is then to determine the distribution of

$$Z = h(X, Y).$$

In Section 3.3 we discussed an approach to this question when the random variables where discrete.

One could prove a result as attained in Exercise 5.3.10 for functions of two variables but this will require knowledge of Linear Algebra and multivariable calculus. Here we limit our objective and shall focus on two specific functions namely the sum and the product.

5.5.1 Distributions of Sums of Independent Random variables

Let X and Y be two independent continuous random variables with densities f_X and f_Y . In this section we shall see how to compute the distribution of $Z = X + Y$. We first prove a proposition that describes the probability density function of Z .

PROPOSITION 5.5.1. (*Sum of two independent random variables*) Let X and Y be two independent random variables with marginal densities given by $f_X : \mathbb{R} \rightarrow \mathbb{R}$ and $f_Y : \mathbb{R} \rightarrow \mathbb{R}$. Then $Z = X + Y$ has a probability density function $f_Z : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x)f_Y(z-x)dx. \quad (5.5.1)$$

Proof- Let us first find an expression for the distribution function of Z .

$$\begin{aligned} F(z) &= P(Z \leq z) \\ &= P(X + Y \leq z) \\ &= \iint_{\{(x,y):x+y \leq z\}} f_X(x)f_Y(y)dydx \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{z-x} f_X(x)f_Y(y)dydx \\ &= \int_{-\infty}^z [\int_{-\infty}^{\infty} f_X(x)f_Y(u-x)dx]du. \end{aligned}$$

As $f_X(\cdot)$ and $f_Y(\cdot)$ are densities, it can be shown that the integrand is a piecewise continuous function. Hence F is of the form (5.2.4) and Theorem 5.2.5 implies that the probability density function of Z is given by (5.5.1). ■

The integral expression on the right hand side of (5.5.1) is referred to as the convolution of f_X and f_Y and is denoted by $f_X \star f_Y(z)$. It is a property of convolutions that $f_X \star f_Y(z) = f_Y \star f_X(z)$ for all $z \in \mathbb{R}$. Thus if we view the sum of X and Y as $Z = X + Y$ or $Z = Y + X$ the distribution will be the same (See Exercise 5.5.8).

EXAMPLE 5.5.2. (*Sum of Uniforms*) Let X and Y be two independent Uniform $(0, 1)$ random variables. Let $Z = X + Y$. From the above proposition that Z has a density given by (5.5.1). Note that

$$f_X(x)f_Y(z-x) = \begin{cases} 1 & \text{if } 0 < x < 1, 0 < z-x < 1 \text{ and } 0 < z < 2 \\ 0 & \text{otherwise} \end{cases}$$

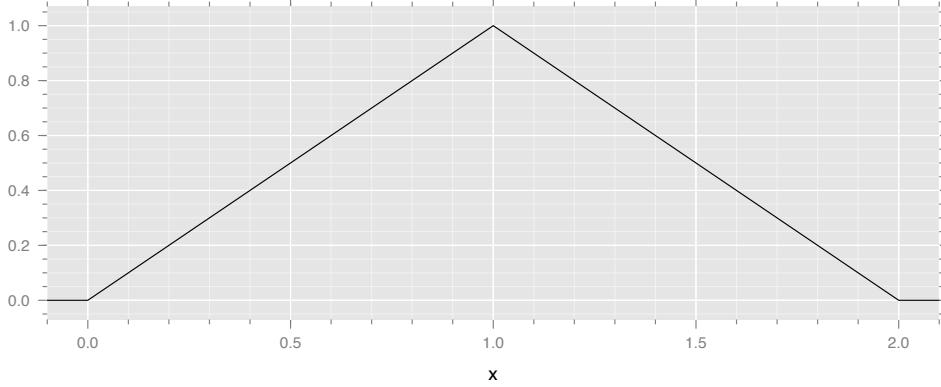
Therefore $f_X(x)f_Y(z-x)$ is non-zero if and only if $\max\{0, z-1\} < x < \min\{1, z\}, 0 < z < 2$. So for $0 < z < 2$,

$$f_Z(z) = \int_{\max\{0, z-1\}}^{\min\{1, z\}} f_X(x)f_Y(z-x)dx = \int_{\max\{0, z-1\}}^{\min\{1, z\}} 1dx = \min\{1, z\} - \max\{0, z-1\}.$$

Therefore,

$$f_Z(z) = \begin{cases} \min\{1, z\} - \max\{0, z-1\} & \text{if } 0 < z < 2 \\ 0 & \text{otherwise} \end{cases} = \begin{cases} z & \text{if } 0 < z \leq 1 \\ 2-z & \text{if } 1 < z < 2 \\ 0 & \text{otherwise} \end{cases}$$

A graph of this density is displayed in Figure 5.14. ■

Figure 5.14: The region $T = \{(x, y) \mid 0 < x < y < 4\}$ from Example 5.5.2.

Our next example will deal with sum of two independent exponential random variables. This will lead us to the Gamma distribution which is of significant interest in statistics.

EXAMPLE 5.5.3. (Sum of Exponentials) Let $\lambda > 0$, X and Y be two independent Exponential (λ) random variables. Let $Z = X + Y$. Then we know and Z has a density given by (5.5.1). Further,

$$f_X(x)f_Y(z-x) = \begin{cases} \lambda^2 e^{-\lambda x} e^{-\lambda(z-x)} & \text{if } x \geq 0, z-x \geq 0 \\ 0 & \text{otherwise} \end{cases} = \begin{cases} \lambda^2 e^{-\lambda z} & \text{if } x \geq 0, x \leq z, z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Hence $f_X(x)f_Y(z-x)$ is non-zero if and only if $0 \leq x \leq z$. So

$$f_Z(z) = \int_0^z f_X(x)f_Y(z-x)dx = \lambda^2 e^{-\lambda z} \int_0^z 1 dx = \lambda^2 z e^{-\lambda z},$$

for $z \geq 0$ and $f_Z(z) = 0$ otherwise. This is known as Gamma $(2, \lambda)$ distribution. ■

Before we define the Gamma distribution more generally we prove a lemma in real analysis, the proof of which can be skipped upon first reading.

LEMMA 5.5.4. For $n \geq 1$, and $\lambda > 0$,

$$\int_0^\infty x^{n-1} e^{-\lambda x} = \frac{(n-1)!}{\lambda^n} \quad (5.5.2)$$

Proof- For all $n \geq 1, \lambda > 0, a > 0$ define $u : [0, a] \rightarrow \mathbb{R}$ and $v : [0, a] \rightarrow \mathbb{R}$ by

$$u(x) = x^{n-1} \text{ and } v(x) = e^{-\lambda x}.$$

As u, v are continuous functions, clearly $I_{n,\lambda}^a$ given by

$$I_{n,\lambda}^a = \int_0^a x^{n-1} e^{-\lambda x}.$$

is well defined finite positive number. As $x^\alpha e^{-\beta x} \rightarrow 0$ as $x \rightarrow \infty$ for any $\alpha, \beta > 0$ there is a $K > 0$ such that

$$0 \leq x^{n-1} e^{-\lambda x} < e^{-\frac{\lambda x}{2}},$$

for all $K > 0$. Therefore $b > a > k$ we have

$$|I_{n,\lambda}^a - I_{n,\lambda}^b| = \int_a^b x^{n-1} e^{-\lambda x} dx \leq \int_a^b e^{-\frac{\lambda x}{2}} dx = 2(e^{-\frac{\lambda b}{2}} - e^{-\frac{\lambda a}{2}}).$$

From this it is standard to note that

$$I_{n,\lambda} := \int_0^\infty x^{n-1} e^{-\lambda x} dx = \lim_{a \rightarrow \infty} I_{n,\lambda}^a$$

is a well defined finite positive number. Now, as u, v are differentiable we have by the integration by parts formula

$$\int_0^a u(x)v'(x)dx = u(a)v(a) - u(0)v(0) - \int_0^a u'(x)v(x)dx.$$

Substituting for u, v above we get

$$-\lambda I_{n,\lambda}^a = a^{n-1} e^{-\lambda a} - (n-1)I_{n-1,\lambda}^a.$$

Taking limits as $a \rightarrow \infty$ we have

$$\lambda I_{n,\lambda} = (n-1)I_{n-1,\lambda}.$$

Applying the above inductively we have

$$I_{n,\lambda} = \prod_{i=1}^{n-1} \frac{(n-i)}{\lambda} I_{1,\lambda} = \frac{(n-1)!}{\lambda^{n-1}} I_{1,\lambda}.$$

Using the fact that $I_1 = \frac{1}{\lambda}$ we have the result. ■

DEFINITION 5.5.5. $X \sim \text{Gamma}(n, \lambda)$: Let $\lambda > 0$ and $n \in \mathbb{N}$. Then X is said to be Gamma distributed with parameters n and λ if it has the density

$$f(x) = \frac{\lambda^n}{(n-1)!} x^{n-1} e^{-\lambda x}, \quad (5.5.3)$$

where $x \geq 0$. The parameter n is referred to as the shape parameter and λ as the rate parameter. By (5.5.2) we know that f given by (5.5.3) is a density function.

We saw in Example 5.5.3 that sum of two exponential distributions resulted in a gamma distribution. If $X \sim \text{Exponential}(\lambda)$ then it can also be viewed as a $\text{Gamma}(1, \lambda)$ distribution. The result in Example 5.5.3 could be rephrased as follows: the sum of two gamma random variables with shape parameter 1 and rate parameter λ is distributed as a gamma random variable with shape parameter 2 and rate parameter λ . This holds more generally as we show in the next example.

EXAMPLE 5.5.6. (Sum of Gammas) Let $n \in \mathbb{N}, m \in \mathbb{N}, \lambda > 0$, X and Y be two independent $\text{Gamma}(n, \lambda)$ and $\text{Gamma}(m, \lambda)$ random variables respectively. Let $Z = X + Y$. Then we know that Z has a density given by (5.5.1). Further,

$$\begin{aligned} f_X(x)f_Y(z-x) &= \begin{cases} \frac{\lambda^n}{(n-1)!} x^{n-1} e^{-\lambda x} \frac{\lambda^m}{(m-1)!} (z-x)^{m-1} e^{-\lambda(z-x)} & \text{if } x \geq 0, z-x \geq 0 \\ 0 & \text{otherwise} \end{cases} \\ &= \begin{cases} \frac{e^{-\lambda z} \lambda^{n+m}}{(n-1)!(m-1)!} x^{n-1} (z-x)^{m-1} & \text{if } x \geq 0, x \leq z, z \geq 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

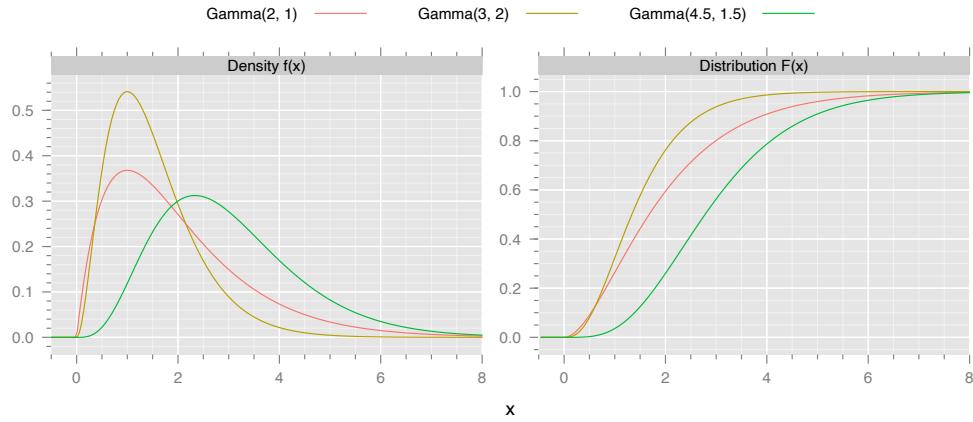


Figure 5.15: The Gamma density and cumulative distribution functions for various shape and rate parameters.

For $z \geq 0$, we have

$$\begin{aligned} f_Z(z) &= \int_{-\infty}^{\infty} f_{X_1}(x)f_{X_2}(z-x)dx = \int_0^z f_{X_1}(x)f_{X_2}(z-x)dx \\ &= \frac{e^{-\lambda z}\lambda^{n+m}}{(n-1)!(m-1)!} \int_0^z x^{n-1}(z-x)^{m-1}dx \end{aligned}$$

We now make a change of variable $x = zu$ so that $dx = zdu$ to obtain

$$f_Z(z) = \frac{z^{n+m-1}e^{-\lambda z}\lambda^{n+m}}{(n-1)!(m-1)!} \int_0^1 u^{n-1}(1-u)^{m-1}du$$

Define

$$c(n, m) = \frac{\int_0^1 u^{n-1}(1-u)^{m-1}du}{(n-1)!(m-1)!}.$$

Thus we have the probability density of Z is given by,

$$f_Z(z) = \begin{cases} c(n, m) \cdot \lambda^{n+m} z^{n+m-1} e^{-\lambda z} & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

To evaluate $c(n, m)$ we use the following fact. From Proposition 5.5.1 $f_Z(\cdot)$ (given by (5.5.1)) is a Probability density function. Therefore,

$$\begin{aligned} 1 &= \int_{-\infty}^{\infty} f_Z(z)dz \\ &= c(n, m) \lambda^{n+m} \int_0^{\infty} z^{n+m-1} e^{-\lambda z} dz \\ &= c(n, m) [(n+m-1)!], \end{aligned}$$

where in the last line we have used (5.5.2) with n replaced by $n+m$. So $c(n, m) = \frac{1}{(n+m-1)!}$. Hence Z has Gamma $(n+m, \lambda)$ distribution. From the definition of $c(n, m)$ we also have

$$\int_0^1 u^{n-1}(1-u)^{m-1}du = \frac{(n+m-1)!}{(n-1)!(m-1)!}.$$

The above calculation is easily extended by an induction argument to obtain the fact that if $\lambda > 0$, X_i , $1 \leq i \leq n$ are independent $\text{Gamma}(n_i, \lambda)$ distributed random variables (respectively). Then $Z = \sum_{i=1}^n X_i$ has $\text{Gamma}(\sum_{i=1}^n n_i, \lambda)$ distribution.

As Exponential (λ) is the same as $\text{Gamma}(1, \lambda)$ random variable, the above implies that the sum of n independent Exponential (λ) random variables is a $\text{Gamma}(n, \lambda)$ random variable. ■

It is possible to define the Gamma distribution when the shape parameter is not necessarily an integer.

DEFINITION 5.5.7. $X \sim \text{Gamma}(\alpha, \lambda)$: Let $\lambda > 0$ and $\alpha > 0$. Then X is said to be Gamma distributed with shape parameter α and rate parameter λ if it has the density

$$f(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, \quad (5.5.4)$$

where $x \geq 0$ and for $\alpha > 0$

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx \quad (5.5.5)$$

One can imitate the calculation done in Example 5.5.6 as well for such a Gamma distribution.

The distribution function of a gamma random variable involves an indefinite form of the integral in (5.5.5). Such integrals are known as incomplete gamma functions, and have no closed-form solution in terms of simple functions. In R, $F(x)$ for the gamma distribution

$$F(x) = P(X \leq x) = \int_0^x \frac{\lambda^\alpha}{\Gamma(\alpha)} z^{\alpha-1} e^{-\lambda z} dz, x > 0$$

can be evaluated numerically with a function call of the form `pgamma(x, alpha, lambda)`. For example,

```
> pgamma(1, 2, 1)
[1] 0.2642411
> pgamma(3, 4.5, 1.5)
[1] 0.5627258
```

Similarly, the density function $f(x)$ in (5.5.4) involves the normalising constant $\Gamma(\alpha)$ (also known as the gamma function) which usually cannot be computed explicitly when α is not an integer. Using R, one can evaluate $f(x)$ numerically using the `dgamma()` function as

```
> dgamma(1, 2, 1)
[1] 0.3678794
> dgamma(3, 4.5, 1.5)
[1] 0.2769272
```

5.5.2 Distributions of Quotients of Independent Random variables.

Let X and Y be two independent continuous random variables with densities f_X and f_Y . In this section we shall find out the probability density function of $Z = \frac{X}{Y}$. As $P(Y = 0) = 0$, Z is well defined random variable.

PROPOSITION 5.5.8. (*Quotient of two independent random variables*) Let X and Y be two independent random variables with marginal densities given by $f_X : \mathbb{R} \rightarrow \mathbb{R}$ and $f_Y : \mathbb{R} \rightarrow \mathbb{R}$. Then $Z = \frac{X}{Y}$ has a probability density function $f_Z : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$f_Z(z) = \int_{-\infty}^{\infty} |y| f_X(zy) f_Y(y) dy. \quad (5.5.6)$$

Proof- Let us find an expression for the distribution function of Z .

$$\begin{aligned} F(z) &= P(Z \leq z) \\ &= P\left(\frac{X}{Y} \leq z\right) \\ &= \int \int_{\{(x,y):y \neq 0, \frac{x}{y} \leq z\}} f_X(x) f_Y(y) dy dx \\ &= \int \int_{\{(x,y):y < 0, \frac{x}{y} \leq z\}} f_X(x) f_Y(y) dy dx + \int \int_{\{(x,y):y > 0, \frac{x}{y} \leq z\}} f_X(x) f_Y(y) dy dx \\ &= \int \int_{\{(x,y):y < 0, x \geq yz\}} f_X(x) f_Y(y) dy dx + \int \int_{\{(x,y):y > 0, x \leq yz\}} f_X(x) f_Y(y) dy dx \\ &= \int_{-\infty}^0 \int_{yz}^{\infty} f_X(x) f_Y(y) dx dy + \int_0^{\infty} \int_{-\infty}^{yz} f_X(x) f_Y(y) dx dy \\ &= I + II \end{aligned}$$

Let us make a u -substitution $x = yu$ in both I and II . For I , $y < 0$, so we will obtain,

$$\begin{aligned} I &= \int_{-\infty}^0 \int_z^{-\infty} y f_X(yu) f_Y(y) du dy \\ &= \int_{-\infty}^0 \int_{-\infty}^z (-y) f_X(yu) f_Y(y) du dy \\ &= \int_{-\infty}^z \int_{-\infty}^0 (-y) f_X(yu) f_Y(y) dy du, \end{aligned}$$

where in the last line we have changed the order of integration¹. For II , $y > 0$ so we will obtain (similarly as in I),

$$\begin{aligned} II &= \int_0^{\infty} \left(\int_{-\infty}^z y f_X(yu) f_Y(y) du \right) dy \\ &= \int_{-\infty}^z \int_0^{\infty} y f_X(yu) f_Y(y) dy du, \end{aligned}$$

Therefore

$$\begin{aligned} F(z) &= I + II \\ &= \int_{-\infty}^z \int_{-\infty}^0 (-y) f_X(yu) f_Y(y) dy du + \int_{-\infty}^z \int_0^{\infty} y f_X(yu) f_Y(y) dy du \\ &= \int_{-\infty}^z \int_{-\infty}^{\infty} |y| f_X(yu) f_Y(y) dy du \end{aligned}$$

¹The change of order of integration is justifiable under certain hypothesis for the integrand. We shall assume these are satisfied, as it is not possible to state and verify them within the scope of this book

As $f_X(\cdot)$ and $f_Y(\cdot)$ are densities, it can be shown that the integrand is a piecewise continuous function. Hence the F is of the form (5.2.4) and Theorem 5.2.5 implies that the probability density function of Z is given by (5.5.6). \blacksquare

Using the above method for finding the distribution of quotient of two random variables, we shall present three examples that will lead us to standard continuous distributions that are useful in applications. We begin with an example that constructs the Cauchy distribution.

EXAMPLE 5.5.9. Let X and Y be two independent Normal random variables with mean 0 and variance $\sigma^2 \neq 0$. Let $Z = \frac{X}{Y}$. We know that the probability density function of Z is given by (5.5.6). Further, for any $y, z \in \mathbb{R}$

$$f_X(zy)f_Y(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{z^2y^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{y^2}{2\sigma^2}} = \frac{1}{2\pi\sigma^2} \exp\left(-\left(\frac{1+z^2}{2\sigma^2}\right)y^2\right)$$

Fix $z \in \mathbb{R}$.

$$\begin{aligned} f_Z(z) &= \int_{-\infty}^{\infty} |y| \frac{1}{2\pi\sigma^2} \exp\left(-\left(\frac{1+z^2}{2\sigma^2}\right)y^2\right) dy \\ &= \frac{1}{2\pi\sigma^2} \left[\int_{-\infty}^0 |y| \exp\left(-\left(\frac{1+z^2}{2\sigma^2}\right)y^2\right) dy + \int_0^{\infty} |y| \exp\left(-\left(\frac{1+z^2}{2\sigma^2}\right)y^2\right) dy \right] \\ &= \frac{1}{2\pi\sigma^2} \left[\int_{-\infty}^0 (-y) \exp\left(-\left(\frac{1+z^2}{2\sigma^2}\right)y^2\right) dy + \int_0^{\infty} y \exp\left(-\left(\frac{1+z^2}{2\sigma^2}\right)y^2\right) dy \right] \end{aligned}$$

It is easy to see that two integrals are the same (perform a substitution of $u = -y$ in the first integral). So the above is

$$= \frac{1}{\pi\sigma^2} \int_0^{\infty} y \exp\left(-\left(\frac{1+z^2}{2\sigma^2}\right)y^2\right) dy.$$

Now perform a substitution $\left(\frac{1+z^2}{2\sigma^2}\right)y^2 = t$, so $\frac{1+z^2}{\sigma^2}ydy = dt$.

$$\begin{aligned} f_Z(z) &= \frac{1}{\pi\sigma^2} \frac{\sigma^2}{1+z^2} \int_0^{\infty} \exp(-t) dt. \\ &= \frac{1}{\pi(1+z^2)} (-e^{-t}) \Big|_0^{\infty} = \frac{1}{\pi(1+z^2)}. \end{aligned}$$

Therefore Z has the Cauchy distribution, which we first saw in the context of Example 5.3.4. \blacksquare

The next example considers the ratio of two gamma random variables. This motivates a standard distribution called the F -distribution, also very important for statistics.

EXAMPLE 5.5.10. Let $m \in \mathbb{N}, n \in \mathbb{N}, \lambda > 0$, X and Y be two independent Gamma (m, λ) and Gamma (n, λ) random variables respectively. Let $Z = \frac{X}{Y}$. We know that the probability density function of Z is given by (5.5.6). Further,

$$\begin{aligned} f_X(zy)f_Y(y) &= \begin{cases} \frac{\lambda^m}{(m-1)!} (zy)^{m-1} e^{-\lambda(zy)} \frac{\lambda^n}{(n-1)!} y^{n-1} e^{-\lambda y} & \text{if } y \geq 0, zy \geq 0 \\ 0 & \text{otherwise} \end{cases} \\ &= \begin{cases} \frac{\lambda^{n+m}}{(n-1)!(m-1)!} y^{n+m-2} z^{m-1} e^{-\lambda(1+zy)} & \text{if } y \geq 0, z \geq 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Fix $z > 0$,

$$\begin{aligned} f_Z(z) &= \int_0^\infty y \frac{\lambda^{n+m}}{(n-1)!(m-1)!} y^{n+m-2} z^{m-1} e^{-\lambda(1+z)y} dy \\ &= \frac{z^{m-1} \lambda^{n+m}}{(n-1)!(m-1)!} \int_0^\infty y^{n+m-1} e^{-\lambda(1+z)y} dy \end{aligned}$$

Now perform a substitution $(1+z)y = t$, so $(1+z)dy = dt$ and the above is

$$= \frac{z^{m-1}}{(1+z)^{m+n}} \frac{\lambda^{n+m}}{(m-1)!(n-1)!} \int_0^\infty t^{m+n-1} e^{-\lambda t} dt$$

Using (5.5.2) we have that

$$f_Z(z) = \begin{cases} \frac{(m+n-1)!}{(m-1)!(n-1)!} z^{m-1} (1+z)^{-(m+n)} & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (5.5.7)$$

■

We shall define the F -distribution in Chapter 8 (See Example 8.1.7) and see further applications of it in Chapter ???. Our next example is a construction of the Beta-distribution.

EXAMPLE 5.5.11. Let $m \in \mathbb{N}, n \in \mathbb{N}, \lambda > 0$. Let X and Y be two independent Gamma (m, λ) and Gamma (n, λ) random variables respectively. Let $Z = \frac{X}{X+Y}$.

Let $W = \frac{Y}{X}$. Note that $Z = \frac{1}{1+W}$. In Example 5.5.10 we found the probability density function of W . We shall use this to find the distribution function of Z . As $P(W \geq 0) = 1$,

$$P(Z \leq z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 1. \end{cases}$$

For $0 < z < 1$,

$$\begin{aligned} P(Z \leq z) &= P\left(\frac{1}{1+W} \leq z\right) = P\left(W \geq \frac{1-z}{z}\right) \\ &= 1 - P\left(W \leq \frac{1-z}{z}\right) \end{aligned}$$

Using (5.5.7) we obtain that the above is

$$\begin{aligned} &= 1 - \int_0^{\frac{1-z}{z}} \frac{(m+n-1)!}{(m-1)!(n-1)!} u^{m-1} (1+u)^{-(m+n)} du \\ &= 1 - \frac{(m+n-1)!}{(m-1)!(n-1)!} \int_0^{\frac{1-z}{z}} u^{m-1} (1+u)^{-(m+n)} du \end{aligned}$$

For $0 < z < 1$, by the fundamental theorem of calculus, differentiating in z

$$\begin{aligned} f_Z(z) &= \frac{1}{z^2} \cdot \frac{(m+n-1)!}{(m-1)!(n-1)!} \left(\frac{1-z}{z}\right)^{m-1} \left(1 + \frac{1-z}{z}\right)^{-(m+n)} \\ &= \frac{(m+n-1)!}{(m-1)!(n-1)!} z^{n-1} (1-z)^{m-1} \end{aligned}$$

Z is said to have the Beta(m, n) distribution.

■

We define the distribution in general next.

DEFINITION 5.5.12. $X \sim \text{Beta}(\alpha, \beta)$: Let $\alpha > 0$ and $\beta > 0$. Then X is said to be Beta distributed with parameters α and β if it has the density

$$f(x) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1} & 0 < x < 1 \\ 0 & \text{otherwise.} \end{cases} \quad (5.5.8)$$

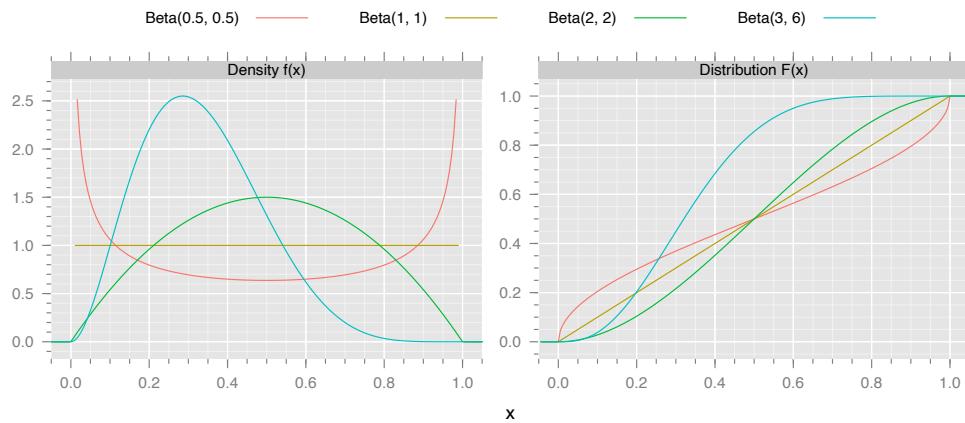


Figure 5.16: The Beta density and cumulative distribution functions for selected shape parameters.

The distribution function of a beta random variable is given by an indefinite integral which in general has no closed-form solution in terms of simple functions. In R, $F(x)$ for the beta distribution

$$F(x) = P(X \leq x) = \int_0^x \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} u^{\alpha-1}(1-u)^{\beta-1} du, \quad 0 < x < 1$$

can be evaluated numerically with a function call of the form `pbeta(x, alpha, beta)`. For example,

```
> pbeta(0.5, 0.5, 0.5)
[1] 0.5
> pbeta(0.5, 3, 6)
[1] 0.8554688
> pbeta(0.2, 6, 1)
[1] 0.2030822
> pbeta(0.2, 1, 6)
[1] 0.737856
```

In the special case where either α or β equals 1, the distribution function of X can be computed explicitly. Another special case is the standard arcsine law we previously encountered in Exercise 5.2.7 in terms of its explicit distribution function; it is easy to see that this is the same as the $\text{Beta}(\frac{1}{2}, \frac{1}{2})$ distribution. The semicircular distribution encountered in Exercise 5.2.6 is also related, in the sense that it can be viewed as a location and scale transformed beta random variable.

EXERCISES

Ex. 5.5.1. Suppose that X and Y are random variables with joint probability density

$$f(x, y) = \begin{cases} \frac{4}{5}(xy + 1) & \text{if } 0 \leq x \leq 1 \text{ and } 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Compute the marginal densities of X and Y ?
- (b) Compute the conditional density $(X|Y = y)$ (for appropriate y).
- (c) Are X and Y independent?

Ex. 5.5.2. Let X and Y be two random variables with the joint p.d.f given by

$$f(x, y) = \begin{cases} \lambda^2 e^{-\lambda y} & 0 \leq x \leq y \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find the marginal distribution of X and Y .
- (b) Find the conditional distribution of $(Y | X = x)$ for some $x > 0$

Ex. 5.5.3. Let $a, b > 0$. Let $X \sim \text{Gamma}(a, b)$ and $Y \sim \text{Exponential}(X)$.

- (a) Find the joint density of X and Y .
- (b) Find the marginal density of Y .
- (c) Find the conditional density of $(X | Y = y)$.

Ex. 5.5.4. Let X_1, X_2, X_3 be independent and identically distributed Uniform $(0, 1)$ random variables. Let $A = X_1 X_3$ and $B = X_2^2$. Find the $P(A < B)$.

Ex. 5.5.5. Let X and Y be two independent exponential random variables each with mean 1.

- (a) Find the density of $U_1 = X^{\frac{1}{2}}$.
- (b) Find the density of $U_2 = X + Y + 1$.
- (c) Find $P(\max\{X, Y\} > 1)$.

Ex. 5.5.6. Suppose X is a uniform random variable in the interval $(0, 1)$ and Y is an independent exponential(2) random variable. Find the distribution of $Z = X + Y$.

Ex. 5.5.7. Let $\alpha > 0, \beta > 0, \lambda > 0$, X and Y be two independent $\text{Gamma}(\alpha, \lambda)$ and $\text{Gamma}(\beta, \lambda)$ random variables respectively. Then $Z = X + Y$ is distributed as a $\text{Gamma}(\alpha + \beta, \lambda)$.

Ex. 5.5.8. Let X and Y be two independent random variables with probability density function $f_X(\cdot)$ and $f_Y(\cdot)$. Show that $X + Y$ and $Y + X$ have the same distribution by showing that the integral expression defining $f_X * f_Y(\cdot)$ is equal to the integral expression defining $f_Y * f_X(\cdot)$.

Ex. 5.5.9. Let $\alpha > 0$ and $\Gamma(\alpha)$ as in (5.5.5).

- (a) Using the same technique as in Lemma 5.5.4, show that $0 < \Gamma(\alpha) < \infty$.
- (b) Show that $\Gamma(\frac{1}{2}) = \int_0^\infty x^{-0.5} e^{-x} dx = \sqrt{\pi}$.

Ex. 5.5.10. Let $\alpha > 0, \delta > 0, \lambda > 0$. Let X and Y be two independent Gamma (α, λ) and Gamma (δ, λ) random variables respectively.

- (a) Let $W = \frac{Y}{X}$. Find the probability density function of W .
- (b) Let $Z = \frac{X}{X+Y}$. Find the probability density function of Z .
- (c) Are X and Z independent ?

Hint: Compute the joint density and see if it is a product of the marginals.

Ex. 5.5.11. Suppose X, Y are independent random variables each normally distributed with mean 0 and variance 1.

- (a) Find the probability density function of $R = \sqrt{X^2 + Y^2}$
- (b) Find the probability density function of $Z = \frac{X}{Y}$
- (c) Find the probability density function of $\theta = \arctan\left(\frac{X}{Y}\right)$
- (d) Are R, θ independent random variables ?

*Hint: Compute the joint density using the change of variable indicated in Exercise 5.1.10.
Decide if it is a product of the marginals*

6

SUMMARISING CONTINUOUS RANDOM VARIABLES

In this chapter we shall revisit concepts that have been discussed for discrete random variables and see their analogues in the continuous setting. We then introduce generating functions and conclude this chapter with a discussion on bivariate normal random variables.

6.1 EXPECTATION, AND VARIANCE

The notion of expected value carries over from discrete to continuous random variables, but instead of being described in terms of sums, it is defined in terms of integrals.

DEFINITION 6.1.1. Let X be a continuous random variable with piecewise continuous density $f(x)$. Then the expected value of X is given by

$$E[X] = \int_{-\infty}^{\infty} xf(x) dx.$$

provided that the integral converges absolutely¹. In this case we say that X has “finite expectation”. If the integral diverges to $\pm\infty$ we say the random variable has infinite expectation. If the integral diverges, but not to $\pm\infty$ we say the expected value is undefined.

The next three examples illustrate the three possibilities: the first is an example where expectation exists as a real number; the next is an example of an infinite expected value; and the final example shows that the expected value may not be defined at all.

EXAMPLE 6.1.2. Let $X \sim \text{Uniform}(a, b)$. Then the expected value of X is given by

$$E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx = \int_a^b x \cdot \frac{1}{b-a} dx = \frac{1}{2(b-a)}(b^2 - a^2) = \frac{b+a}{2}.$$

This result is intuitive since it says that the average value of a $\text{Uniform}(a, b)$ random variable is the midpoint of its interval. ■

EXAMPLE 6.1.3. Let $0 < \alpha < 1$ and $X \sim \text{Pareto}(\alpha)$ which is defined to have the probability density function

$$f(x) = \begin{cases} \frac{\alpha}{x^{\alpha+1}} & 1 \leq x < \infty \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = \int_1^{\infty} x \cdot \frac{\alpha}{x^{\alpha+1}} dx = \alpha \lim_{M \rightarrow \infty} \int_1^M x^{-\alpha} dx = \frac{\alpha}{-\alpha+1} (-1 + \lim_{M \rightarrow \infty} M^{-\alpha+1}) = \infty$$

as $0 < \alpha < 1$.

Thus this Pareto random variable has an infinite expected value. ■

EXAMPLE 6.1.4. Let $X \sim \text{Cauchy}(0, 1)$. Then the probability density function of X is given by

$$f(x) = \frac{1}{\pi} \frac{1}{1+x^2} \text{ for all } x \in \mathbb{R}.$$

Now,

$$E[X] = \int_{-\infty}^{\infty} x \cdot \frac{1}{\pi(1+x^2)} dx$$

Now by Exercise 6.1.10, we know that as $M \rightarrow -\infty, N \rightarrow \infty$ the $\int_M^N \frac{x}{1+x^2} dx$ does not converge or diverge to $\pm\infty$. So $E[X]$ is not defined for this Cauchy random variable. ■

Expected values of functions of continuous random variables may be computed using their respective probability density function by the following theorem.

THEOREM 6.1.5. Let X be continuous random variables with probability density function $f_X : \mathbb{R} \rightarrow \mathbb{R}$.

(a) Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be piecewise continuous and $Z = g(X)$. Then the expected value of Z given by

$$E[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x) dx$$

(b) Let Y be a continuous random variable such that (X, Y) have a joint probability density function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Suppose $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ be piecewise continuous. Then,

$$E[h(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} h(x, y)f(x, y) dx dy.$$

Proof- The proof is beyond the scope of this book. For (a) when g is as in Exercise 5.3.10 then one can provide the proof using only the tools of basic calculus (we will leave this case as an exercise to the reader) ■

We illustrate the use of the above theorem with a couple of examples.

EXAMPLE 6.1.6. A piece of equipment breaks down after a functional lifetime that is a random variable $T \sim \text{Exp}(\frac{1}{5})$. An insurance policy purchased on the equipment pays a dollar amount equal to $1000 - 200t$ if the equipment breaks down at a time $0 \leq t \leq 5$ and pays nothing if the equipment breaks down after time $t = 5$. What is the expected payment of the insurance policy?

For $t \geq 0$ the policy pays $g(t) = \max\{1000 - 200t, 0\}$ so,

$$\begin{aligned} E[g(T)] &= \int_0^{\infty} \frac{1}{5} e^{-(1/5)t} \max\{1000 - 200t, 0\} dt \\ &= \int_0^5 \frac{1}{5} e^{-(1/5)t} (1000 - 200t) dt \\ &= 1000e^{-1} \approx \$367.88 \end{aligned}$$

EXAMPLE 6.1.7. Let $X, Y \sim \text{Uniform}(0, 1)$. What is the expected value of the larger of the two variables?

We offer two methods of solving this problem. The first is to define $Z = \max\{X, Y\}$ and then determine the density of Z . To do so, we first find its distribution. $F_Z(z) = P(Z \leq z)$, but

$\max\{X, Y\}$ is less than or equal to z exactly when both X and Y are less than or equal to z . So for $0 \leq z \leq 1$,

$$\begin{aligned} F_Z(z) &= P((X \leq z) \cap (Y \leq z)) \\ &= P(X \leq z) \cdot P(Y \leq z) \\ &= z^2 \end{aligned}$$

Therefore $f_Z(z) = F'_Z(z) = 2z$ after which the expected value can be obtained through integration

$$E[Z] = \int_0^1 z \cdot 2z \, dz = \frac{2}{3} z^3 \Big|_0^1 = \frac{2}{3}.$$

An alternative method is to use Theorem 6.1.5 (b) to calculate the expectation directly without finding a new density. Since X and Y are independent, their joint distribution is the product of their marginal distributions. That is,

$$f(x, y) = f_X(x)f_Y(y) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1 \text{ and } 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Therefore,

$$\begin{aligned} E[\max\{X, Y\}] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \max\{x, y\} \cdot f(x, y) \, dx \, dy \\ &= \int_0^1 \int_0^1 \max\{x, y\} \cdot 1 \, dx \, dy \end{aligned}$$

The value of $\max\{x, y\}$ is x if $0 < y \leq x < 1$ and it is y if $0 < x \leq y < 1$. So,

$$\begin{aligned} E[\max\{X, Y\}] &= \int_0^1 \int_0^y y \, dx \, dy + \int_0^1 \int_y^1 x \, dx \, dy \\ &= \int_0^1 xy \Big|_{x=0}^{x=y} \, dy + \int_0^1 \frac{1}{2}x^2 \Big|_{x=y}^{x=1} \, dy \\ &= \int_0^1 y^2 \, dy + \int_0^1 \frac{1}{2} - \frac{1}{2}y^2 \, dy \\ &= \frac{1}{3} + \frac{1}{3} = \frac{2}{3}. \end{aligned}$$

■

Results from calculus may be used to show that the linearity properties from Theorem 4.1.7 such as apply to continuous random variables as well as to discrete ones. We restate it here for completeness.

THEOREM 6.1.8. *Suppose that X and Y are continuous random variables with piecewise continuous joint density function function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Assume that both have finite expected value. If a and b are real numbers then*

- (a) $E[aX] = aE[X]$;
- (b) $E[aX + b] = aE[X] + b$
- (c) $E[X + Y] = E[X] + E[Y]$; and

$$(d) E[aX + bY] = aE[X] + bE[Y].$$

$$(e) \text{ If } X \geq 0 \text{ then } E[X] \geq 0.$$

Proof- See Exercise 6.1.11. ■

We will use these now-familiar properties in the continuous setting. As in the discrete setting we can define the variance and standard deviation of a continuous random variable.

DEFINITION 6.1.9. Let X be a random variable with probability density function $f : \mathbb{R} \rightarrow \mathbb{R}$. Suppose X has finite expectation. Then

(a) the variance of the random variable is written as $\text{Var}[X]$ and is defined as

$$\text{Var}[X] = E[(X - E[X])^2] = \int_{-\infty}^{\infty} (x - E[X])^2 f_X(x) dx,$$

(b) the standard deviation of X is written as $SD[X]$ and is defined as

$$SD[X] = \sqrt{\text{Var}[X]}$$

Since the above terms are expected values, there is the possibility that they may be infinite because the integral describing the expectation diverges to infinity. As the integrand is strictly positive, it isn't possible for the integral to diverge unless it diverges to infinity.

The properties of variance and standard deviation of continuous random variables match those of their discrete counterparts. A list of these properties follows below.

THEOREM 6.1.10. Let $a \in \mathbb{R}$ and let X be a continuous random variable with finite variance (and thus, with finite expected value as well). Then,

$$(a) \text{Var}[X] = E[X^2] - (E[X])^2.$$

$$(b) \text{Var}[aX] = a^2 \cdot \text{Var}[X];$$

$$(c) SD[aX] = |a| \cdot SD[X];$$

$$(d) \text{Var}[X + a] = \text{Var}[X]; \text{ and}$$

$$(e) SD[X + a] = SD[X].$$

If Y is another independent continuous random variable with finite variance (and thus, with finite expected value as well) then

$$(f) E[XY] = E[X]E[Y];$$

$$(g) \text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]; \text{ and}$$

$$(h) SD[X + Y] = \sqrt{(\text{SD}[X])^2 + (\text{SD}[Y])^2}.$$

Proof- The proof is essentially an imitation of the proofs presented in Theorem 4.1.10, Theorem 4.2.5, Theorem 4.2.4, and Theorem 4.2.6. One needs to use the respective densities, integrals in lieu of sums, and use Theorem 6.1.11 and Theorem 6.1.5 when needed. We will leave this as an exercise to the reader. ■

EXAMPLE 6.1.11. Let $X \sim \text{Normal}(0, 1)$. In this example we shall show that $E[X] = 0$ and $\text{Var}[X] = 1$. Before that we collect some facts about the probability density function of X , given by (5.2.7). Using (5.2.9) with $z = 0$, we can conclude that

$$\int_0^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \frac{1}{2} \quad (6.1.1)$$

Observe that there exists $c_1 > 0$ such that

$$\max\{|x|, x^2\} e^{-\frac{x^2}{2}} \leq c_1 e^{-c_1|x|}$$

for all $x \in \mathbb{R}$. Hence

$$\begin{aligned} \int_{-\infty}^\infty |x| \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx &\leq c_1 \int_{-\infty}^\infty e^{-c_1|x|} dx < \infty \\ \int_{-\infty}^\infty x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx &\leq c_1 \int_{-\infty}^\infty e^{-c_1|x|} dx < \infty \end{aligned} \quad (6.1.2)$$

Using the above we see that

$$E[X] = \int_{-\infty}^\infty x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx < \infty$$

So we can split integral expression in definition of $E[X]$ as

$$E[X] = \int_{-\infty}^0 x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx + \int_0^\infty x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx.$$

Further the change of variable $y = -x$ will imply that

$$\int_{-\infty}^0 x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = - \int_0^\infty y \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy.$$

So $E[X] = 0$. Again by (6.1.2),

$$\text{Var}[X] = \int_{-\infty}^\infty (x - E[X])^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_{-\infty}^\infty x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx < \infty$$

To evaluate the integral we make a change of variable to obtain

$$\int_{-\infty}^\infty x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_{-\infty}^0 x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx + \int_0^\infty x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = 2 \int_0^\infty x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx.$$

Then we use integration by parts like Lemma 5.5.4. Set $u(x) = x$ and $v(x) = e^{-\frac{x^2}{2}}$, which imply $u'(x) = 1$ and $v'(x) = -xe^{-\frac{x^2}{2}}$. Therefore for $a > 0$,

$$\begin{aligned} \int_0^a x^2 e^{-\frac{x^2}{2}} dx &= \int_0^a u(x)(-v'(x)) dx = u(x)(-v(x)) \Big|_0^a - \int_0^a u'(x)(-v(x)) dx \\ &= a^2 e^{-\frac{a^2}{2}} + \int_0^a e^{-\frac{x^2}{2}} \end{aligned}$$

Using the fact that $\lim_{a \rightarrow \infty} a^2 e^{-\frac{a^2}{2}} = 0$ and (6.1.1) we have

$$\begin{aligned} \text{Var}[X] &= 2 \frac{1}{\sqrt{2\pi}} \int_0^\infty x^2 e^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{\pi}} \lim_{a \rightarrow \infty} \int_0^a x^2 e^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{\pi}} \lim_{a \rightarrow \infty} \left[a^2 e^{-\frac{a^2}{2}} + \int_0^a e^{-\frac{x^2}{2}} dx \right] \\ &= \frac{1}{\sqrt{\pi}} \left[0 + \int_0^\infty e^{-\frac{x^2}{2}} dx \right] = \frac{1}{\sqrt{\pi}} [0 + \sqrt{\pi}] = 1 \end{aligned}$$

Suppose $Y \sim \text{Normal}(\mu, \sigma^2)$ then we know by Corollary 5.3.3 that $W = \frac{Y-\mu}{\sigma} \sim \text{Normal}(0, 1)$. By Example 6.1.11, $E[W] = 0$ and $\text{Var}[W] = 1$. Also $Y = \sigma W + \mu$, so by Theorem 6.1.8(b) $E[Y] = \sigma E[W] + \mu = \mu$ and by Theorem 6.1.10 (d) and (b) $\text{Var}[Y] = \sigma^2 \text{Var}[W] = \sigma^2$. \blacksquare

EXAMPLE 6.1.12. Let $X \sim \text{Uniform}(a, b)$. To calculate the variance of X first note that Theorem 6.1.5(a) gives

$$E[X^2] = \int_{-\infty}^{\infty} x^2 \cdot f(x) dx = \int_a^b x^2 \cdot \frac{1}{b-a} dx = \frac{1}{3(b-a)}(b^3 - a^3) = \frac{b^2 + ab + a^2}{3}.$$

Now, since $E[X] = \frac{b+a}{2}$ (see Example 6.1.2), the variance may be found as

$$\text{Var}[X] = E[X^2] - (E[X])^2 = \frac{b^2 + ab + a^2}{3} - \left(\frac{b+a}{2}\right)^2 = \frac{(b-a)^2}{12}.$$

Taking square roots, we obtain $SD[X] = \frac{b-a}{\sqrt{12}}$. So the standard deviation of a continuous, uniform random variable is $\frac{1}{\sqrt{12}}$ times of the length of its interval. \blacksquare

The Markov and Chebychev inequalities also apply to continuous random variables. As with discrete variables, these help to estimate the probabilities that a random variable will fall within a certain number of standard deviations from its expected value.

THEOREM 6.1.13. *Let X be a continuous random variable with probability density function f and finite non-zero variance.*

(a) **(Markov's Inequality)** Suppose X is supported on non-negative values, i.e. $f(x) = 0$ for all $x < 0$. Then for any $c > 0$,

$$P(X \geq c) \leq \frac{\mu}{c}.$$

(b) **(Chebychev's Inequality)** For any $k > 0$,

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

Proof - (a) By definition of μ and assumptions on f , we have

$$\mu = \int_{-\infty}^{\infty} xf(x)dx = \int_0^{\infty} xf(x)dx.$$

Using an elementary fact from integrals we know that

$$\int_0^{\infty} xf(x)dx = \int_0^c xf(x)dx + \int_c^{\infty} xf(x)dx$$

We note that the first integral is non-negative so we have

$$\mu \geq \int_c^{\infty} xf(x)dx.$$

As $f(\cdot) \geq 0$, we have $xf(x) \geq cf(x)$ whenever $x > c$. So again using facts about integrals

$$\mu \geq \int_c^{\infty} cf(x)dx = c \int_c^{\infty} f(x)dx = cP(X > c).$$

The last equality follows from definition. Hence we have the result.

(b) The event ($|X - \mu| \geq k\sigma$) is the same as the event ($(X - \mu)^2 \geq k^2\sigma^2$). The random variable $(X - \mu)^2$ is certainly non-negative, is continuous by Exercise 5.3.9, and its expected value is the variance of X which we have assumed to be finite. Therefore we may apply Markov's inequality to $(X - \mu)^2$ to get

$$P(|X - \mu| \geq k\sigma) = P((X - \mu)^2 \geq k^2\sigma^2) \leq \frac{E[(X - \mu)^2]}{k^2\sigma^2} = \frac{\text{Var}[X]}{k^2\sigma^2} = \frac{\sigma^2}{k^2\sigma^2} = \frac{1}{k^2}.$$

Though the theorem is true for all $k > 0$, it doesn't give any useful information unless $k > 1$. ■

EXERCISES

Ex. 6.1.1. Suppose X has probability density function given by

$$f_X(x) = \begin{cases} 1 - |x| & -1 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Compute the distribution function of X .
- (b) Compute $E[X]$ and $\text{Var}[X]$.

Ex. 6.1.2. Suppose X has probability density function given by

$$f_X(x) = \begin{cases} \frac{\cos(x)}{2} & -\frac{\pi}{2} \leq x \leq \frac{\pi}{2} \\ 0 & \text{otherwise} \end{cases}$$

- (a) Compute the distribution function of X .
- (b) Compute $E[X]$ and $\text{Var}[X]$.

Ex. 6.1.3. Find $E[X]$ and $\text{Var}[X]$ in the following situations:

- (a) $X \sim \text{Normal}(\mu, \sigma^2)$, with $\mu \in \mathbb{R}$ and $\sigma > 0$.
- (b) X has probability density function given by

$$f_X(x) = \begin{cases} x & 0 \leq x \leq 1 \\ 2 - x & 1 \leq x \leq 2 \\ 0 & \text{otherwise} \end{cases}$$

Ex. 6.1.4. Let $1 < \alpha$ and $X \sim \text{Pareto}(\alpha)$. Calculate $E[X]$ to show that it is finite.

Ex. 6.1.5. Let X be a random variable with density $f(x) = 2x$ for $0 < x < 1$ (and $f(x) = 0$ otherwise).

- (a) Calculate $E[X]$. You should get a result larger than $\frac{1}{2}$. Explain why this should be expected even without computations.
- (b) Calculate $SD[X]$.

Ex. 6.1.6. Let $X \sim \text{Uniform}(a, b)$ and let $k > 0$. Let μ and σ be the expected value and standard deviation calculated in Example 6.1.12.

- (a) Calculate $P(|X - \mu| \leq k\sigma)$. Your final answer should depend on k , but not on the values of a or b .
- (b) What is the value of k such that results of more than k standard deviations from expected value are unachievable for X ?

Ex. 6.1.7. Let $X \sim \text{Exponential}(\lambda)$.

- (a) Prove that $E[X] = \frac{1}{\lambda}$ and $SD[X] = \frac{1}{\lambda}$.
- (b) Let μ and σ denote the mean and standard deviation of X respectively. Use your computations from (a) to calculate $P(|X - \mu| \leq k\sigma)$. Your final answer should depend on k , but not on the value of λ .
- (c) Is there a value of k such that results of more than k standard deviations from expected value are unachievable for X ?

Ex. 6.1.8. Let $X \sim \text{Gamma}(n, \lambda)$ with $n \in \mathbb{N}$ and $\lambda > 0$. Using Example 5.5.3, Exercise 6.1.7(a) and Theorem 6.1.8(c) calculate $E[X]$. Using Theorem 6.1.10 calculate $Var[X]$.

Ex. 6.1.9. Let $X \sim \text{Uniform}(0, 10)$ and let $g(x) = \max\{x, 4\}$. Calculate $E[g(X)]$.

Ex. 6.1.10. Show that as $M \rightarrow -\infty, N \rightarrow \infty$ $\int_M^N \frac{x}{1+x^2} dx$ does not have a limit.

Ex. 6.1.11. Using the hints provided below prove the respective parts of Theorem 6.1.8.

- (a) For $a = 0$ the result is clear. Let $a \neq 0$ and $f_X : \mathbb{R} \rightarrow \mathbb{R}$ be the probability density function of X . Use Lemma 5.3.2 to find the probability density function of aX . Compute the expectation of aX to obtain the result. Alternatively use Theorem 6.1.5(a).
- (b) Use Theorem 6.1.5(b).
- (c) Use the joint density of (X, Y) to write $E[X + Y]$. Then use (5.4.2) and (5.4.3) to prove the result.
- (d) Use the same technique as in (b).
- (e) If $X \geq 0$ then its marginal density $f_X : \mathbb{R} \rightarrow \mathbb{R}$ is positive only when the $x \geq 0$. The result immediately follows from definition of expectation.

Ex. 6.1.12. Prove Theorem 6.1.10.

6.2 COVARIANCE, CORRELATION, CONDITIONAL EXPECTATION AND CONDITIONAL VARIANCE

Covariance of continuous random variables (X, Y) is used to describe how the two random variables relate to each other. The properties proved about covariances for discrete random variables in Section 4.5 apply to continuous random variables as well via essentially the same arguments. We define covariance and state the properties next.

DEFINITION 6.2.1. Let X and Y be random variables with joint probability density function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Suppose X and Y have finite expectation. Then the covariance of X and Y is defined as

$$\text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E[X])(y - E[Y])f(x, y)dxdy, \quad (6.2.1)$$

Since it is defined in terms of an expected value, there is the possibility that the covariance may be infinite or not defined at all. We now state the properties of Covariance.

THEOREM 6.2.2. *Let X, Y be continuous random variables such that they have joint probability density function. Assume that $0 \neq \sigma_x^2 = \text{Var}(X) < \infty, 0 \neq \sigma_y^2 = \text{Var}(Y) < \infty$. Then*

- (a) $\text{Cov}[X, Y] = E[XY] - E[X]E[Y]$.
- (b) $\text{Cov}[X, Y] = \text{Cov}[Y, X]$;
- (c) $\text{Cov}[X, X] = \text{Var}[X]$.
- (d) $-\sigma_X\sigma_Y \leq \text{Cov}[X, Y] \leq \sigma_X\sigma_Y$
- (e) If X and Y are independent then $\text{Cov}[X, Y] = 0$.

Let a, b be real numbers. Suppose Z is another continuous random variable, and $\sigma_z = \text{Var}(Z) < \infty$. Further (X, Z) , (Y, Z) , $(X, aY + bZ)$, and $(aX + bY, Z)$ all have (their respective) joint probability functions. Then

- (f) $\text{Cov}[X, aY + bZ] = a \cdot \text{Cov}[X, Y] + b \cdot \text{Cov}[X, Z]$;
- (g) $\text{Cov}[aX + bY, Z] = a \cdot \text{Cov}[X, Z] + b \cdot \text{Cov}[Y, Z]$;

Proof- See Exercise 6.2.13 ■

DEFINITION 6.2.3. Let (X, Y) be continuous random variables both with finite variance and covariance. From Theorem 6.2.2(d) the quantity $\rho[X, Y] = \frac{\text{Cov}[X, Y]}{\sigma_X\sigma_Y}$ is in the interval $[-1, 1]$. It is known as the “correlation” of X and Y . As discussed earlier, both the numerator and denominator include the units of X and the units of Y . The correlation, therefore, has no units associated with it. It is thus a dimensionless rescaling of the covariance and is frequently used as an absolute measure of trends between the two continuous random variables as well.

EXAMPLE 6.2.4. Let $X \sim \text{Uniform}(0, 1)$ and be independent of $Y \sim \text{Uniform}(0, 1)$. Let $U = \min(X, Y)$ and $V = \max(X, Y)$. We wish to find $\rho[U, V]$. First, $0 < u < 1$

$$P(U \leq u) = 1 - P(U > u) = 1 - P(X > u, Y > u) = 1 - P(X > u)P(Y > u) = 1 - (1-u)^2,$$

as X, Y are independent uniform random variables. Second, for $0 < v < 1$,

$$P(V \leq v) = P(X \leq v, Y \leq v) = P(X \leq v)P(Y \leq v) = v^2,$$

as X, Y are independent uniform random variables. Therefore the distribution function of U and V are given by

$$F_U(u) = \begin{cases} 0 & \text{if } u < 0 \\ 1 - (1-u)^2 & \text{if } 0 < u < 1 \\ 1 & \text{if } u \geq 1. \end{cases} \quad \text{and} \quad F_V(v) = \begin{cases} 0 & \text{if } v < 0 \\ v^2 & \text{if } 0 < v < 1 \\ 1 & \text{if } v \geq 1. \end{cases}$$

As F_U, F_V are piecewise differentiable, the probability density function of U and V are obtained by differentiating F_U and F_V respectively.

$$f_U(u) = \begin{cases} 2(1-u) & \text{if } 0 < u < 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad f_V(v) = \begin{cases} 2v & \text{if } 0 < v < 1 \\ 0 & \text{otherwise.} \end{cases}$$

Thirdly, $0 < u < v < 1$

$$\begin{aligned} P(U \leq u, V \leq v) &= P(V \leq v) - P(U > u, V \leq v) \\ &= v^2 - P(u < X \leq v, u < Y \leq v) \\ &= v^2 - P(u < X \leq v)P(u < Y \leq v) \\ &= v^2 - (v-u)^2, \end{aligned}$$

where we have used the formula for distribution function of V and the fact that X, Y are independent uniform random variables. It is easily seen that $P(U \leq u, V \leq v) = 0$ for all other possibilities of (u, v) . As the joint distribution function is piecewise differentiable in each variable, the joint probability density function of U and V , $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, exists and is obtained by differentiating it partially in u and v .

$$f(u, v) = \begin{cases} 2 & \text{if } 0 < u < v < 1 \\ 0 & \text{otherwise} \end{cases}$$

Now,

$$\begin{aligned} E[U] &= \int_0^1 u 2(1-u) du = u^2 - 2\frac{u^3}{3} \Big|_0^1 = \frac{1}{3} \\ E[V] &= \int_0^1 v 2v dv = 2\frac{v^3}{3} \Big|_0^1 = \frac{2}{3} \\ E[U^2] &= \int_0^1 u^2 2(1-u) du = 2\frac{u^3}{3} - 2\frac{u^4}{4} \Big|_0^1 = \frac{1}{6} \\ E[V^2] &= \int_0^1 v^2 2v dv = 2\frac{v^4}{4} \Big|_0^1 = \frac{1}{2} \\ E[UV] &= \int_0^1 \left[\int_0^v uv 2du \right] dv = \int_0^1 2v \left[\frac{u^2}{2} \Big|_0^1 \right] dv = \int_0^1 2v \frac{v^2}{2} dv = \frac{v^4}{4} \Big|_0^1 = \frac{1}{4} \end{aligned}$$

Therefore

$$\begin{aligned} Var[U] &= E[U^2] - (E[U])^2 = \frac{2}{3} - \frac{1}{9} = \frac{5}{9} \\ Var[V] &= E[V^2] - (E[V])^2 = \frac{1}{2} - \frac{4}{9} = \frac{1}{18} \\ Cov[U, V] &= E[UV] - E[U]E[V] = \frac{1}{4} - \frac{1}{3}\frac{2}{3} = \frac{5}{36} \\ \rho[U, V] &= \frac{Cov[U, V]}{\sqrt{Var[V]}\sqrt{Var[U]}} = \frac{\frac{5}{36}}{\sqrt{\frac{5}{9}}\sqrt{\frac{1}{18}}} = \frac{1}{2\sqrt{2}} \end{aligned}$$

■

As seen in Theorem 6.2.2 (e), independence of X and Y guarantees that they are uncorrelated (i.e $\rho[X, Y] = 0$). The converse is not true (See Example 4.5.6 for discrete case). It is possible that $Cov[X, Y] = 0$ and yet that X and Y are dependent, as the next example shows.

EXAMPLE 6.2.5. Let $X \sim \text{Uniform}(-1, 1)$. Let $Y = X^2$. Note from Example 6.1.2 and Example 6.1.12 we have $E[X] = 0, E[Y] = E[X^2] = \frac{1}{3}$. Further using the probability density function of X ,

$$E[XY] = E[X^3] = \int_{-1}^1 x^3 \frac{1}{2} dx = \frac{x^4}{8} \Big|_{-1}^1 = 0.$$

So $\rho[X, Y] = 0$. Clearly X and Y are not independent. We verify this precisely as well. Consider the

$$P(X \leq -\frac{1}{4}, Y \leq \frac{1}{4}) = P(X \leq -\frac{1}{4}, X^2 \leq \frac{1}{4}) = P(-\frac{1}{2} \leq X \leq -\frac{1}{4}) = \frac{1}{8},$$

as $X \sim \text{Uniform } (-1, 1)$. Whereas,

$$P(X \leq -\frac{1}{4})P(Y \leq \frac{1}{4}) = P(X \leq -\frac{1}{4})P(X^2 \leq \frac{1}{4}) = P(X \leq -\frac{1}{4})P(-\frac{1}{2} \leq X \leq \frac{1}{2}) = \frac{3}{8} \cdot \frac{1}{2} = \frac{3}{16}.$$

Clearly

$$P(X \leq -\frac{1}{4}, Y \leq \frac{1}{4}) \neq P(X \leq -\frac{1}{4})P(Y \leq \frac{1}{4})$$

implying they are not independent. ■

We are now ready to define conditional expectation and variance.

DEFINITION 6.2.6. Let (X, Y) be continuous random variables with a piecewise continuous joint probability density function f . Let f_X be the marginal density of X . Assume x is a real number for which $f_X(x) \neq 0$. The conditional expectation of Y given $X = x$ is defined by

$$E[Y | X = x] = \int_{-\infty}^{\infty} y f_{Y|X=x}(y) dy = \int_{-\infty}^{\infty} y \frac{f(x, y)}{f_X(x)} dy$$

whenever it exists. The conditional variance of Y given $X = x$ is defined by

$$\begin{aligned} \text{Var}[Y | X = x] &= E[(Y - E[Y | X = x])^2 | X = x] \\ &= \int_{-\infty}^{\infty} \left(y - \int_{-\infty}^{\infty} y \frac{f(x, y)}{f_X(x)} dy \right)^2 \frac{f(x, y)}{f_X(x)} dy. \end{aligned}$$

The results proved in Theorem 4.4.4, Theorem 4.4.6, Theorem 4.4.8, and Theorem 4.4.9 are all applicable when X and Y are continuous random variables having joint probability density function f . The proofs of these results in the continuous setting follow very similarly (though using facts about integrals from analysis).

THEOREM 6.2.7. Let (X, Y) be continuous random variables with joint probability density function $f : \mathbb{R} \rightarrow \mathbb{R}$. Assume that $h, g : \mathbb{R} \rightarrow \mathbb{R}$ be defined as

$$g(y) = \begin{cases} E[X | Y = y] & \text{if } f_Y(y) > 0 \\ 0 & \text{otherwise} \end{cases} \quad \text{and } h(y) = \begin{cases} \text{Var}[X | Y = y] & \text{if } f_Y(y) > 0 \\ 0 & \text{otherwise} \end{cases}$$

are well-defined piecewise continuous functions. Let $k : \mathbb{R} \rightarrow \mathbb{R}$ be a piecewise continuous function. Then

$$E[k(X) | Y = y] = \int_{-\infty}^{\infty} k(x) f_{X|Y=y}(x) dx, \tag{6.2.2}$$

$$E[g(Y)] = E[X], \tag{6.2.3}$$

and

$$\text{Var}[X] = E[h(Y)] + \text{Var}[g(Y)]. \tag{6.2.4}$$

Proof- The proof of (6.2.2) is beyond the scope of this book. We shall omit it. To prove (6.2.3) we use the definition of g and Theorem 6.1.8 (a) to write

$$E[g(Y)] = \int_{-\infty}^{\infty} g(y) f_Y(y) dy = \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} x f_{X|Y=y}(x) dx \right] f_Y(y) dy$$

Using the definition of conditional density and rearranging the order of integration we obtain that the above is

$$= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} x \frac{f(x,y)}{f_Y(y)} dx \right] f_Y(y) dy = \int_{-\infty}^{\infty} x \left[\int_{-\infty}^{\infty} f(x,y) dy \right] dx = \int_{-\infty}^{\infty} x f_X(x) dx = E[X].$$

So we are done. To prove (6.2.4), using Exercise 6.2.8

$$h(y) = E[X^2 | Y = y] - (E[X | Y = y])^2 = E[X^2 | Y = y] - (g(y))^2$$

From the above we have,

$$\begin{aligned} E[h(Y)] &= E[X^2] + E[g(Y)^2] \\ Var[g(Y)] &= E[g(Y)^2] - (E[g(Y)])^2 = E[g(Y)^2] - (E[X])^2 \end{aligned}$$

Therefore summing the two equations we have (6.2.4). ■

As before it is common to use $E[X|Y]$ to denote $g(Y)$ after which the result may be expressed as $E[E[X|Y]] = E[X]$. This can be slightly confusing notation, but one must keep in mind that the exterior expected value in the expression $E[E[X|Y]]$ refers to the average of $E[X|Y]$ viewed as a function of Y .

Similarly one denotes $h(Y)$ by $Var[X|Y]$. Then we can rewrite (6.2.4) as

$$Var[X] = E[Var[X|Y]] + Var[E[X|Y]].$$

EXAMPLE 6.2.8. Let $X \sim \text{Uniform}(0, 1)$ and be independent of $Y \sim \text{Uniform}(0, 1)$. Let $U = \min(X, Y)$ and $V = \max(X, Y)$. In Example 6.2.4 we found $\rho[U, V]$. During that computation we showed that the marginal densities of U and V were given by

$$f_U(u) = \begin{cases} 2(1-u) & \text{if } 0 < u < 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad f_V(v) = \begin{cases} 2v & \text{if } 0 < v < 1 \\ 0 & \text{otherwise.} \end{cases}$$

and the joint density of (U, V) was given by

$$f(u, v) = \begin{cases} 2 & \text{if } 0 < u < v < 1 \\ 0 & \text{otherwise} \end{cases}$$

Let $0 < u < 1$. The conditional density of $V | U = u$, is given by

$$f_{V|U=u}(v) = \frac{f(u, v)}{f_U(u)}, \text{ for } v \in \mathbb{R}.$$

So,

$$f_{V|U=u}(v) = \begin{cases} \frac{1}{1-u} & \text{if } u < v < 1 \\ 0 & \text{otherwise} \end{cases}$$

Therefore $(V | U = u) \sim \text{Uniform}(u, 1)$. So the conditional expectation is given by

$$E[V | U = u] = \int_u^1 \frac{v}{1-u} dv = \frac{1-u^2}{2(1-u)} = \frac{1+u}{2}.$$

The conditional variance is given by

$$\begin{aligned} \text{Var}[V | U = u] &= E[V^2 | U = u] - (E[V | U = u])^2 \\ &= \int_u^1 \frac{v^2}{1-u} dv - \left(\frac{1+u}{2}\right)^2 \\ &= \frac{1-u^3}{3(1-u)} dv - \frac{(1+u)^2}{4} = \frac{(1-u)^2}{12}. \end{aligned}$$

We could have also concluded these from properties of Uniform distribution computed in Example 6.1.2 and Example 6.1.12. We will use this approach in the next example. ■

EXAMPLE 6.2.9. Let (X, Y) have joint probability density function f given by

$$f(x, y) = \frac{\sqrt{3}}{4\pi} e^{-\frac{1}{2}(x^2-xy+y^2)} \quad -\infty < x, y < \infty.$$

These random variables were considered in Example 5.4.12. We showed there that X is a Normal random variable with mean 0 and variance $\frac{4}{3}$ and Y is also a Normal random variable with mean 0 and variance $\frac{4}{3}$. We observed that they are not independent as well and the conditional distribution of Y given $X = x$ was Normal with mean $\frac{x}{2}$ and variance 1. Either by direct computation or by definition we observe that

$$E[Y | X = x] = \frac{x}{2} \quad \text{Var}[Y | X = x] = 1.$$

We could compute the $\text{Var}[Y]$ using (6.2.4), i.e.

$$\begin{aligned} \text{Var}[Y] &= \text{Var}[E[Y | X]] + E[\text{Var}[Y | X = x]] \\ &= \text{Var}\left[\frac{X}{2}\right] + E[1] \\ &= \frac{1}{4}\text{Var}[X] + 1 = \frac{1}{4}\frac{4}{3} + 1 = \frac{4}{3}. \end{aligned}$$

■

EXERCISES

Ex. 6.2.1. Let (X, Y) be uniformly distributed on the triangle $0 < x < y < 1$.

- (a) Compute $E[X|Y = \frac{1}{6}]$.
- (b) Compute $E[(X - Y)^2]$.

Ex. 6.2.2. X is a random variable with mean 3 and variance 2. Y is a random variable with mean -1 and variance 6. The covariance of X and Y is -2 . Let $U = X + Y$ and $V = X - Y$. Find the correlation coefficient of U and V .

Ex. 6.2.3. Suppose X and Y are both uniformly distributed on $[0, 1]$. Suppose $\text{Cov}[X, Y] = \frac{-1}{24}$. Compute the variance of $X + Y$.

Ex. 6.2.4. A dice game between two people is played by a pair of dice being thrown. One of the dice is green and the other is white. If the green die is larger than the white die, player number one earns a number of points equal to the value on the green die. If the green die is less than or equal to the white die, then player number two earns a number of points equal to the value of the green die. Let X be the random variable representing the number of points earned by player one after one throw. Let Y be the random variable representing the number of points earned by player two after one throw.

- (a) Compute the expected value of X and of Y .
- (b) Without explicitly computing it, would you expect $\text{Cov}[X, Y]$ to be positive or negative? Explain.
- (c) Calculate $\text{Cov}[X, Y]$ to confirm your intuition.

Ex. 6.2.5. Suppose X has variance σ_X^2 , Y has variance σ_Y^2 , and the pair (X, Y) has correlation coefficient $\rho[X, Y]$.

- (a) In terms of σ_X , σ_Y , and $\rho[X, Y]$, find $\text{Cov}[X, Y]$ and $\text{Cov}[X + Y, X - Y]$.
- (b) What must be true of σ_X^2 and σ_Y^2 if $X + Y$ and $X - Y$ are uncorrelated?

Ex. 6.2.6. Let (X, Y) have the joint probability density function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by

$$f_{X,Y}(x, y) = \begin{cases} 3(x+y) & \text{if } x > 0, y > 0, \text{ and } x+y < 1 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find $E[X|Y = \frac{1}{2}]$ and $\text{Var}[X|Y = \frac{1}{2}]$
- (b) Are X and Y independent?

Ex. 6.2.7. Suppose Y is uniformly distributed on $(0, 1)$, and suppose for $0 < y < 1$ the conditional density of $X | Y = y$ is given by

$$f_{X|Y=y}(x) = \begin{cases} \frac{2x}{y^2} & \text{if } 0 < x < y \\ 0 & \text{otherwise.} \end{cases}$$

- (a) Show that, as a function of x , $f_{X|Y=y}(x)$ is a density.
- (b) Compute the joint p.d.f. of (X, Y) and the marginal density of X .
- (c) Compute the expected value and variance of X given that $Y = y$, with $0 < y < 1$.

Ex. 6.2.8. Let (X, Y) have joint probability density function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Show that $\text{Var}[X | Y = y] = E[X^2 | Y = y] - (E[X | Y = y])^2$.

Ex. 6.2.9. For random variables (X, Y) as in Exercise 5.4.1, find

- (a) $E[X]$ and $E[Y]$
- (b) $\text{Var}[X]$ and $\text{Var}[Y]$
- (c) $\text{Cov}[X, Y]$ and $\rho[X, Y]$

Ex. 6.2.10. From Example 5.4.12, consider (X, Y) have joint probability density function f given by

$$f(x, y) = \frac{\sqrt{3}}{4\pi} e^{-\frac{1}{2}(x^2 - xy + y^2)} \quad -\infty < x, y < \infty.$$

Find

- (a) $E[X]$ and $E[Y]$
- (b) $\text{Var}[X]$ and $\text{Var}[Y]$

- (c) $\text{Cov}[X, Y]$ and $\rho[X, Y]$

Ex. 6.2.11. From Example 5.4.13, suppose $T = \{(x, y) \mid 0 < x < y < 4\}$ and let $(X, Y) \sim \text{Uniform}(T)$. Find

- (a) $E[X]$ and $E[Y]$
- (b) $\text{Var}[X]$ and $\text{Var}[Y]$
- (c) $\text{Cov}[X, Y]$ and $\rho[X, Y]$

Ex. 6.2.12. From Example 5.4.9, consider the open disk in \mathbb{R}^2 given by $C = \{(x, y) : x^2 + y^2 < 25\}$ and $|C| = 25\pi$ denote its area. Let (X, Y) have a joint density $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by

$$f(x, y) = \begin{cases} \frac{1}{|C|} & \text{if } (x, y) \in C \\ 0 & \text{otherwise.} \end{cases}$$

Find

- (a) $E[X]$ and $E[Y]$
- (b) $\text{Var}[X]$ and $\text{Var}[Y]$
- (c) $\text{Cov}[X, Y]$ and $\rho[X, Y]$

Ex. 6.2.13. Using the hints provided below prove the respective parts of Theorem 6.2.2

- (a) Use the linearity properties of the expected value from Theorem 6.1.8.
- (b) Use definition of covariance.
- (c) Use the definitions of variance and covariance.
- (d) Imitate the proof of Theorem 4.5.7.
- (e) Use part (a) of this problem and part (f) of Theorem ??.
- (f) Use the linearity properties of the expected value from Theorem 6.1.8.
- (g) Use the linearity properties of the expected value from Theorem 6.1.8.

Ex. 6.2.14. Let X, Y be continuous random variable with piecewise continuous densities $f(x)$ and $g(y)$ and well-defined expected values. Suppose $X \leq Y$ then show that $E[X] \leq E[Y]$.

Ex. 6.2.15. Let T be the triangle bounded by the lines $y = 0$, $y = 1 - x$, and $y = 1 + x$. Suppose a random vector (X, Y) has a joint p.d.f.

$$f_{(X,Y)}(x, y) = \begin{cases} 3y & \text{if} \\ 0 & \text{otherwise.} \end{cases}$$

Compute $E[Y|X = \frac{1}{2}]$.

Ex. 6.2.16. Let (X, Y) be random variables with joint probability density function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Assume that both random variables have finite variances and that their covariance is also finite.

- (a) Show that $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$.
- (b) Show that when X and Y are positively correlated (i.e. $\rho[X, Y] > 0$) then $\text{Var}[X + Y] > \text{Var}[X] + \text{Var}[Y]$, while when X and Y are negatively correlated (i.e. $\rho[X, Y] < 0$), then $\text{Var}[X + Y] < \text{Var}[X] + \text{Var}[Y]$.

6.3 MOMENT GENERATING FUNCTIONS

We have already seen for the distribution of a discrete random variable or a continuous random variable is determined by its distribution function. In this section we shall discuss the concept of moment generating functions. Under suitable assumptions, these functions will determine the distribution of random variables. They are also serve as tools in computations and come in handy for convergence concepts that we will discuss.

The moment generating function generates or determine the moments which in turn, under suitable hypothesis determine the distribution of the corresponding random variable. We begin with a definition of a moment.

DEFINITION 6.3.1. Suppose X is a random variable. For a positive integer k , the quantity

$$m_k = E[X^k]$$

is known as the " k^{th} moment of X ". As before the existence of a given moment is determined by whether the above expectation exists or not.

We have previously seen many computations of the first moment $E[X]$ and also seen that the second moment $E[X^2]$ is related to the variance of the random variable. The next theorem states that if a moment exists then it guarantees the existence of all lesser moments.

THEOREM 6.3.2. Let X be a random variable and let k be a positive integer. If $E[X^k] < \infty$ then $E[X^j] < \infty$ for all positive integers $j < k$.

Proof - Suppose X is a continuous random variable. Suppose $E[X^k]$ exists and is finite, so that $E[|X^k|] < \infty$. Divide \mathbb{R} in two pieces by letting $R_1 = \{x \in T : |x| < 1\}$ and letting $R_2 = \{x \in T : |x| \geq 1\}$. If $j < k$ then $|x|^j \leq |x|^k$ for $x \in R_2$ so,

$$\begin{aligned} E[|X^j|] &= \int_{\mathbb{R}} |x|^j f_X(x) dx = \int_{R_1} |x|^j f_X(x) dx + \int_{R_2} |x|^j f_X(x) dx \\ &\leq \int_{R_1} 1 \cdot f_X(x) dx + \int_{R_2} |x|^k f_X(x) dx \\ &\leq \int_{R_1} f_X(x) dx + \int_{R_2} |x|^k f_X(x) dx \\ &= 1 + E[|X^k|] < \infty \end{aligned}$$

Therefore $E[X^j]$ exists and is finite. See Exercise 6.3.7 when X is a discrete random variable. ■

When a random variable has finite moments for all positive integers, then these moments provide a great deal of information about the random variable itself. In fact, in some cases, these moments serve to completely describe the distribution of the random variable. One way to simultaneously describe all moments of such a variable in terms of a single expression is through the use of a "moment generating function".

DEFINITION 6.3.3. Suppose X is a random variable and $D = \{t \in \mathbb{R} : E[e^{tX}] \text{ exists}\}$. The function $M : D \rightarrow \mathbb{R}$ given by

$$M(t) = E[e^{tX}],$$

is called the moment generating function for X .

The notation $M_X(t)$ will also be used when clarification is needed as to which variable a particular moment generating function belongs. Note that $M(0) = 1$ will always be true, but for other values of t , there is no guarantee that the function is even defined as the expected value might be infinite. However, when $M(t)$ has derivatives defined at zero, these values incorporate information about the moments of X . For a discrete random variable $X : S \rightarrow T$ with $T = \{x_i : i \in \mathbb{N}\}$, then for $t \in D$ (as in Definition 6.3.3)

$$M_X(t) = \sum_{i \geq 1} e^{tx_i} P(X = x_i).$$

For a continuous random variable X with probability density function $f_X : \mathbb{R} \rightarrow \mathbb{R}$ then for $t \in D$ (as in Definition 6.3.3)

$$M_X(t) = \int_{\mathbb{R}} e^{tx} f_X(x) dx.$$

We compute moment generating function for a Poisson (λ) and a Gamma (n, λ), with $n \in \mathbb{N}, \lambda > 0$.

EXAMPLE 6.3.4. Suppose $X \sim \text{Poisson } (\lambda)$ then for all $t \in \mathbb{R}$,

$$M_X(t) = \sum_{k=0}^{\infty} e^{tk} P(X = k) = \sum_{k=0}^{\infty} e^{tk} \frac{\lambda^k e^{-\lambda}}{k!} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{(e^t \lambda)^k}{k!} = e^{-\lambda} e^{e^t \lambda} = e^{-\lambda(1+e^t)}.$$

So the moment generating function of X exists for all $t \in \mathbb{R}$. Suppose $Y \sim \text{Gamma } (n, \lambda)$ then $t < \lambda$,

$$M_Y(t) = \int_R e^{ty} \frac{\lambda^n}{\Gamma(n)} y^{n-1} e^{-\lambda y} dy = \frac{\lambda^n}{\Gamma(n)} \int_R y^{n-1} e^{-(\lambda-t)y} dy = \frac{\lambda^n}{\Gamma(n)} \frac{\Gamma(n)}{(\lambda-t)^n} = \left(\frac{\lambda}{\lambda-t} \right)^n,$$

where we have used (5.5.3). The moment generating function of Y will not be finite if $t \geq \lambda$. ■

We summarily compile some facts about moment generating functions. The proof of some of the results are beyond the scope of this text.

THEOREM 6.3.5. Suppose for a random variable X , there exists $\delta > 0$ such that $M_X(t)$ exists $(-\delta, \delta)$.

(a) The k^{th} moment of X exists and is given by

$$E[X^k] = M_X^{(k)}(0),$$

where $M_X^{(k)}$ denotes the k^{th} derivative of M_X .

(b) For $0 \neq a \in \mathbb{R}$ such that $a, t \in (-\delta, \delta)$ we have

$$M_{aX}(t) = M_X(at).$$

(c) Suppose Y is another independent random variable such that $M_Y(t)$ exists for $t \in (-\delta, \delta)$. Then

$$M_{X+Y}(t) = M_X(t)M_Y(t).$$

for $t \in (-\delta, \delta)$.

Proof - (a) A precise proof is beyond the scope of this book. We provide a sketch. Express e^{tX} as a power series in t .

$$e^{tX} = 1 + tX + \frac{t^2 X^2}{2} + \cdots + \frac{t^n X^n}{n!} + \cdots$$

The expected value of the left hand side is the moment generating function for X while linearity may be used on the right hand side. So the power series of $M(t)$ is given by

$$M(t) = 1 + t \cdot E[X] + \frac{t^2}{2} \cdot E[X^2] + \cdots + \frac{t^n}{n!} \cdot E[X^n] + \cdots$$

Taking k derivatives of both sides of the equation (which is valid in the interval of convergence) yields

$$M^{(k)}(t) = E[X^k] + t \cdot E[X^{k+1}] + \frac{t^2}{2} \cdot E[X^{k+2}] + \cdots$$

Finally, when evaluating both sides at $t = 0$ all but one term on the right hand side vanishes and the equation becomes simply $M^{(k)}(0) = E[X^k]$.

(b) $M_{aX}(t) = E[e^{(aX)t}] = E[e^{X(at)}] = M_X(at)$.

(c) Using Theorem 4.1.10 or Theorem 6.1.10 (f) we have

$$M_{X+Y}(t) = E[e^{t(X+Y)}] = E[e^{tX}e^{tY}] = E[e^{tX}]E[e^{tY}] = M_X(t)M_Y(t).$$

■

Theorem 6.3.5 applies equally well for both discrete and continuous variables. A discrete example is presented next.

EXAMPLE 6.3.6. Let $X \sim \text{Geometric}(p)$. We shall find $M_X(t)$ and use this function to calculate the expected value and variance of X . For any $t \in \mathbb{R}$,

$$\begin{aligned} M_X(t) &= E[e^{tX}] = \sum_{n=1}^{\infty} e^{tn} P(X = n) = \sum_{n=1}^{\infty} (e^t)^n \cdot p(1-p)^{n-1} = pe^t \cdot \sum_{n=1}^{\infty} (e^t \cdot (1-p))^{n-1} \\ &= \frac{pe^t}{1 - e^t(1-p)} \end{aligned}$$

Having completed that computation, the expected value and variance can be computed simply by calculating derivatives.

$$M'_X(t) = \frac{pe^t}{[1 - (1-p)e^t]^2}$$

and so $E[X] = M'_X(0) = \frac{p}{p^2} = \frac{1}{p}$. Similarly,

$$M''_X(t) = \frac{pe^t + p(1-p)e^{2t}}{[1 - (1-p)e^t]^3}$$

and so $E[X^2] = M''_X(0) = \frac{2p-p^2}{p^3} = \frac{2}{p^2} - \frac{1}{p}$. Therefore, $\text{Var}[X] = E[X^2] - (E[X])^2 = \frac{1-p}{p^2}$. Both the expected value and variance are in agreement with the previous computations for the geometric random variable.

Let $Y \sim \text{Normal}(\mu, \sigma^2)$. The density of Y is $f_Y(y) = \frac{1}{\sigma\sqrt{2\pi}}e^{-(y-\mu)^2/2\sigma^2}$. For any $t \in \mathbb{R}$,

$$\begin{aligned} M_Y(t) &= E[e^{tY}] = \int_{-\infty}^{\infty} e^{ty} \cdot \frac{1}{\sigma\sqrt{2\pi}}e^{-(y-\mu)^2/2\sigma^2} dy = \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}}e^{-(y^2-(2\mu y+2\sigma^2 ty)+\mu^2)/2\sigma^2} dy \\ &= e^{\mu t+(1/2)\sigma^2 t^2} \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}}e^{-(y-(\mu+\sigma^2 t))^2/2\sigma^2} dy \\ &= e^{\mu t+(1/2)\sigma^2 t^2} \end{aligned} \tag{6.3.1}$$

where the integral in the final step is equal to one since it integrates the density of a $\text{Normal}(\mu + \sigma^2 t, \sigma^2)$ random variable. One can easily verify that the $M'_Y(0) = \mu$ and $M''_Y(0) = \mu^2 + \sigma^2$. ■

As with the expected value and variance, moment generating functions behave well when applied to linear combinations of independent variables (courtesy Theorem 6.3.5 (b) and (c)).

EXAMPLE 6.3.7. Suppose we wish to find the moment generating function of $X \sim \text{Binomial}(n, p)$. We have seen that such a random variable may arise as the sum of independent Bernoulli variables. That is, $X = Y_1 + \dots + Y_n$ where $Y_j \sim \text{Bernoulli}(p)$. But it is routine to compute

$$M_{Y_j}(t) = E[e^{tY_j}] = e^{t \cdot 1} P(Y_j = 1) + e^{t \cdot 0} P(Y_j = 0) = pe^t + (1 - p).$$

Therefore by linearity (inductively applying Theorem 6.3.5 (c)),

$$M_X(t) = M_{Y_1 + \dots + Y_n}(t) = M_{Y_1}(t) \cdot \dots \cdot M_{Y_n}(t) = (pe^t + (1 - p))^n.$$

■

Moment generating functions are an extraordinarily useful tool in analyzing the distributions of random variables. Two particularly useful tools involve the uniqueness and limit properties of such generating functions. Unfortunately these theorems require analysis beyond the scope of this text to prove. We will state the uniqueness fact (unproven) below and the limit property in Chapter 8. First we generalize the definition of moment generating functions to pairs of random variables.

DEFINITION 6.3.8. Suppose X and Y are random variables. Then the function

$$M(s, t) = E[e^{sX+tY}]$$

is called the (joint) moment generating function for X and Y . The notation $M_{X,Y}(s, t)$ will be used when confusion may arise as to which random variables are being represented.

Moment generating functions completely describe the distributions of random variables. We state the result precisely.

THEOREM 6.3.9. (M.G.F. Uniqueness Theorem)

- (a) (One variable) Suppose X and Y are random variables and $M_X(t) = M_Y(t)$ in some open interval containing the origin. Then X and Y are equal in distribution.
- (b) (Two variable) Suppose (X, W) and (Y, Z) are pairs of random variables and suppose $M_{X,W}(s, t) = M_{Y,Z}(s, t)$ in some rectangle containing the origin. Then (X, W) and (Y, Z) have the same joint distribution.

An immediate application of the theorem is an alternate proof of Corollary 5.3.3 based on moment generating functions.

EXAMPLE 6.3.10. Let $X \sim \text{Normal}(\mu, \sigma^2)$ and let $Y = \frac{X-\mu}{\sigma}$. Show that $Y \sim \text{Normal}(0, 1)$.

We know X is normal, (6.3.1) shows that the moment generating function of X is $M_X(t) = e^{\mu t + (1/2)\sigma^2 t^2}$, for all $t \in \mathbb{R}$. So consider the moment generating function of Y . For all $t \in \mathbb{R}$

$$\begin{aligned} M_Y(t) &= E[e^{tY}] = E[e^{t(X-\mu)/\sigma}] = E[e^{tX/\sigma} e^{-t\mu/\sigma}] = e^{-t\mu/\sigma} \cdot M_X\left(\frac{t}{\sigma}\right) \\ &= e^{-t\mu/\sigma} \cdot e^{\mu(t/\sigma) + (1/2)\sigma^2(t/\sigma)^2} = e^{\frac{t^2}{2}}. \end{aligned}$$

But this expression is the moment generating function of a $\text{Normal}(0, 1)$ random variable. So by the uniqueness of moment generating functions, Theorem 6.3.9 (a), the distribution of Y is $\text{Normal}(0, 1)$. ■

Just as the joint density of a pair of random variables factors as a product of marginal densities exactly when the variables are independent (Theorem 5.4.7), a similar result holds for moment generating functions.

THEOREM 6.3.11. *Suppose (X, Y) are a pair of continuous random variables with moment generating function $M(s, t)$. Then X and Y are independent if and only if*

$$M(s, t) = M_X(s) \cdot M_Y(t).$$

Proof - One direction of the proof follows from basic facts about independence. If X and Y are independent, then by Exercise 6.3.4, we have

$$M(s, t) = E[e^{sX+tY}] = E[e^{sX}e^{tY}] = E[e^{sX}]E[e^{tY}] = M_X(s) \cdot M_Y(t).$$

To prove the opposite direction, we shall use Theorem 6.3.9(b). Let \hat{X} and \hat{Y} be independent, but have the same distributions as X and Y respectively. Since $M_{X,Y}(s, t) = M_X(s)M_Y(t)$ we have the following series of equalities:

$$M_{X,Y}(s, t) = M_X(s)M_Y(t) = M_{\hat{X}}(s)M_{\hat{Y}}(t) = M_{\hat{X}, \hat{Y}}(s, t).$$

By Theorem 6.3.9(b), this means that (X, Y) and (\hat{X}, \hat{Y}) have the same distribution. This would imply that

$$P(X \in A, Y \in B) = P(\hat{X} \in A, \hat{Y} \in B) = P(\hat{X} \in A)P(\hat{Y} \in B) = P(X \in A)P(Y \in B),$$

for any events A and B . Hence X and Y are independent. ■

Notice that the method employed in Example 6.3.10 did not require considering integrals directly. Since the manipulation of integrals can be complicated (particularly when dealing with multiple integrals), the moment generating function method will often be simpler as the next example illustrates.

EXAMPLE 6.3.12. Let a, b be two real numbers. Let $X \sim \text{Normal}(\mu_1, \sigma_1^2)$ and $Y \sim \text{Normal}(\mu_2, \sigma_2^2)$ be independent. Observe that

$$M_{aX+bY}(t) = M_{X,Y}(at, bt)$$

Using Theorem 6.3.11, we have that the above is

$$M_X(at)M_Y(bt) = e^{a\mu_1 t + (1/2)a^2\sigma_1^2 t^2} e^{b\mu_2 t + (1/2)b^2\sigma_2^2 t^2} = e^{(a\mu_1 + b\mu_2)t + (1/2)(a^2\sigma_1^2 + b^2\sigma_2^2)t^2}$$

which is the moment generating function of a Normal random variable with mean $a\mu_1 + b\mu_2$ and variance $a^2\sigma_1^2 + b^2\sigma_2^2$. So $aX + bY \sim \text{Normal}(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2)$. ■

We conclude this section with a result on finite linear combinations of independent normal random variables.

THEOREM 6.3.13. *Let X_1, X_2, \dots, X_n be independent, normally distributed random variables with mean μ_i and variance σ_i^2 respectively for $i = 1, 2, \dots, n$. Let a_1, a_2, \dots, a_n be real-valued numbers, not all of which are zero. Then the linear combination $Y = a_1X_1 + a_2X_2 + \dots + a_nX_n$ is also normally distributed with mean $\sum_{i=1}^n a_i\mu_i$ and variance $\sum_{i=1}^n a_i^2\sigma_i^2$.*

Proof- This follows from the preceding example by induction and is left as an exercise. ■

EXERCISES

Ex. 6.3.1. Let $X \sim \text{Normal}(0, 1)$. Use the moment generating function of X to calculate $E[X^4]$.

Ex. 6.3.2. Let $Y \sim \text{Exponential}(\lambda)$.

- (a) Calculate the moment generating function $M_Y(t)$.
- (b) Use (a) to calculate $E[Y^3]$ and $E[Y^4]$, the third and fourth moments of an exponential distribution.

Ex. 6.3.3. Let X_1, X_2, \dots, X_n be i.i.d. random variables.

- (a) Let $Y = X_1 + \dots + X_n$. Prove that $M_Y(t) = [M_{X_1}(t)]^n$.
- (b) Let $Z = (X_1 + \dots + X_n)/n$. Prove that $M_Z(t) = [M_{X_1}(\frac{t}{n})]^n$.

Ex. 6.3.4. Let X and Y be two independent discrete random variables. Let $h : \mathbb{R} \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$. Show that

$$E[h(X)g(Y)] = E[h(X)]E[g(Y)].$$

Show that the above holds if X and Y are independent continuous random variables.

Ex. 6.3.5. Suppose X is a discrete random variable and $D = \{t \in \mathbb{R} : E[t^X] \text{ exists}\}$. The function $\psi : D \rightarrow \mathbb{R}$ given by

$$\psi(t) = E[t^X],$$

is called the probability generating function for X . Calculate the probability generating function of X when X is

- (a) $X \sim \text{Bernoulli}(p)$, with $0 < p < 1$.
- (b) $X \sim \text{Binomial}(n, p)$, with $0 < p < 1$, $n \geq 1$.
- (c) $X \sim \text{Geometric}(p)$, with $0 < p < 1$.
- (d) $X \sim \text{Poisson } (\lambda)$, with $0 < \lambda$.

Ex. 6.3.6. Let $X, Y : S \rightarrow T$ be discrete random variables with the number of elements in T is finite. Prove part (a) of Theorem 6.3.9 in this case.

Ex. 6.3.7. Prove Theorem 6.3.2 when X is a discrete random variable.

6.4 BIVARIATE NORMALS

In Example 6.3.12, we saw that if X and Y are independent, normally distributed random variables, any linear combination $aX + bY$ is also normally distributed. In such a case the joint density of (X, Y) is determined easily (courtesy Theorem 5.4.7). We would like to understand random variables that are not independent but have normally distributed marginals. Motivated by the observations in Example 6.3.12 we provide the following definition.

DEFINITION 6.4.1. A pair of random variables (X, Y) is called “bivariate normal” if $aX + bY$ is a normally distributed random variable for all real numbers a and b .

We need to be somewhat cautious in the above definition. Since the variables are dependent it may turn out that $aX + bY = 0$ or some constant. (E.g: $Y = -X$, or $Y = -X + 2$ with $a = 1, b = 1$). We shall follow the convention that a constant c random variable in such cases is a normal random variable with mean c and variance 0.

If (X, Y) are bivariate normal then as $X = X + 0Y$ and $Y = 0X + Y$ both X and Y individually are normal random variables. The converse if not true (See Exercise 6.4.3). However the joint distribution of bivariate normal random variables are determined by their means, variances and covariances. This fact is proved next.

THEOREM 6.4.2. Suppose (X, Y) and (Z, W) are two bivariate normal random variables. If

$$\begin{aligned} E[X] &= E[Z] = \mu_1, & E[Y] &= E[W] = \mu_2 \\ Var[X] &= Var[Z] = \sigma_1^2, & Var[Y] &= Var[W] = \sigma_2^2 \\ && \text{and} \\ Cov[X, Y] &= Cov[Z, W] = \sigma_{12} \end{aligned} \tag{6.4.1}$$

then (X, Y) and (Z, W) have the same joint distribution.

Proof- As (X, Y) and (Z, W) are bivariate normal random variables, given real numbers s, t $sX + tY$ and $sZ + tW$ are normal random variables. Using (6.4.1) and the properties of mean and covariance (see Theorem 6.2.2) we have

$$\begin{aligned} E[sX + tY] &= sE[X] + tE[Y] = s\mu_1 + t\mu_2, \\ E[sZ + tW] &= sE[Z] + tE[W] = s\mu_1 + t\mu_2, \\ Var[sX + tY] &= s^2Var[X] + t^2Var[Y] + 2stCov[X, Y] \\ &= s^2\sigma_1^2 + t^2\sigma_2^2 + 2st\sigma_{12}, \\ &\text{and} \\ Var[sZ + tW] &= s^2Var[Z] + t^2Var[W] + 2stCov[Z, W] \\ &= s^2\sigma_1^2 + t^2\sigma_2^2 + 2st\sigma_{12}. \end{aligned}$$

From the above, $sX + tY$ and $sZ + tW$ have the same mean and variance. So they have the same distribution (as normal random variables are determined by their mean and variances). By Theorem 6.3.9 (a) they have the same moment generating function. So, the (joint) moment generating function of (X, Y) at (s, t) is

$$M_{X,Y}(s, t) = E[e^{sX+tY}] = M_{sX+tY}(1) = M_{sZ+tW}(1) = E[e^{sZ+tW}] = M_{Z,W}(s, t)$$

Therefore (Z, W) has the same joint m.g.f. as (X, Y) and Theorem 6.3.9 (b) implies that they have the same joint distribution. ■

Though, in general, two variables which are uncorrelated may not be independent, it is a remarkable fact that the two concepts are equivalent for bivariate normal random variables.

THEOREM 6.4.3. Let (X, Y) be a bivariate normal random variable. Then $Cov[X, Y] = 0$ if and only if X and Y are independent.

Proof - That independence implies a zero covariance is true for any pair of random variables (use Theorem 6.1.10 (e)), so we need to only consider the reverse implication.

Suppose $Cov[X, Y] = 0$. Let μ_X and σ_X^2 denote the expected value and variance of X and μ_Y and σ_Y^2 the corresponding values for Y . Let s and t be real numbers. Then, by the bivariate

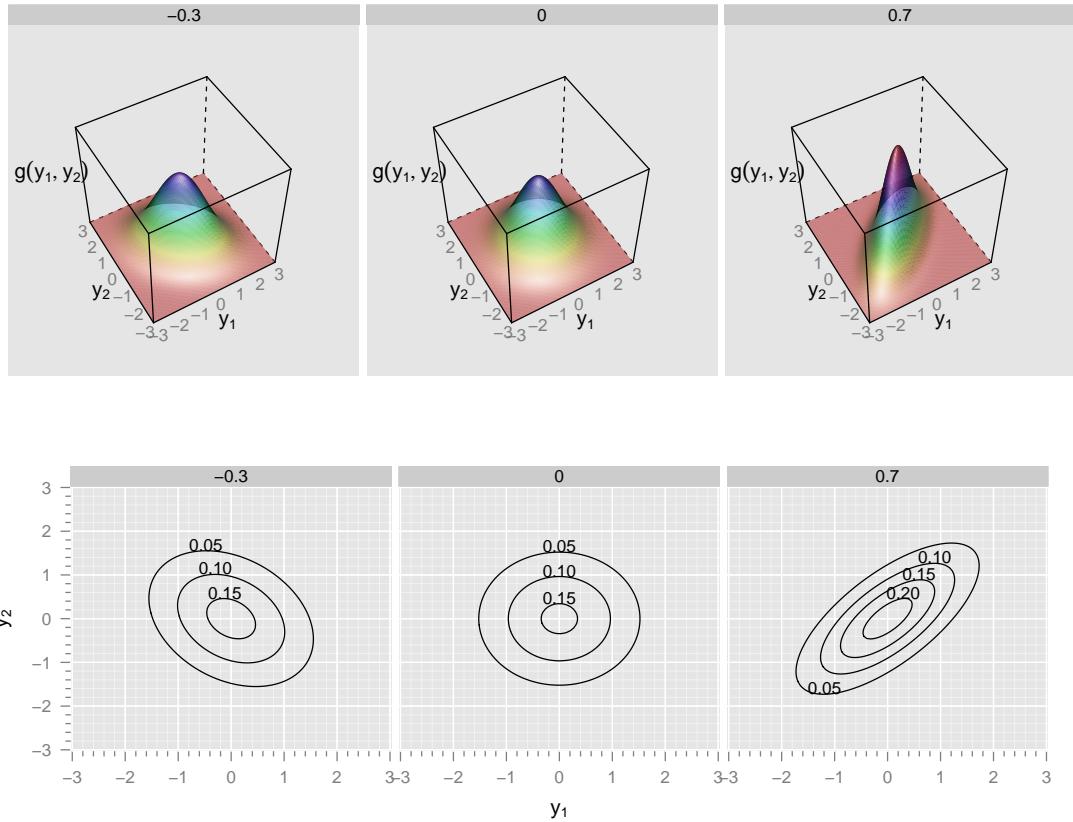


Figure 6.1: The density function of Bivariate Normal distributions. The set of panels on top show a three-dimensional view of the density function for various values of the correlation ρ . The bottom set of panels show contour plots, where each ellipse corresponds to the (y_1, y_2) pairs corresponding to a constant value of $g(y_1, y_2)$.

normality of (X, Y) , we know $sX + tY$ is normally distributed. Moreover by properties of expected value and variance we have

$$E[sX + tY] = sE[X] + tE[Y] = s\mu_X + t\mu_Y$$

and

$$\text{Var}[sX + tY] = s^2\text{Var}[X] + 2st\text{Cov}[X, Y] + t^2\text{Var}[Y] = s^2\sigma_X^2 + t^2\sigma_Y^2.$$

That is, $sX + tY \sim \text{Normal}(s\mu_X + t\mu_Y, s^2\sigma_X^2 + t^2\sigma_Y^2)$. So for all $s, t \in \mathbb{R}$

$$\begin{aligned} M_{X,Y}(s, t) &= E[e^{sX+tY}] = M_{sX+tY}(1) = e^{(s\mu_X+t\mu_Y)+(1/2)(s^2\sigma_X^2+t^2\sigma_Y^2)} \\ &= e^{s\mu_X+(1/2)s^2\sigma_X^2} \cdot e^{t\mu_Y+(1/2)t^2\sigma_Y^2} \\ &= M_X(s) \cdot M_Y(t). \end{aligned}$$

Hence by Theorem 6.3.11 X and Y are independent. ■

We conclude this section by finding the joint density of a Bivariate normal random variable. See Figure 6.1 for a graphical display of this density.

THEOREM 6.4.4. Let (Y_1, Y_2) be a bivariate Normal random variable, with $\mu_1 = E[Y_1], \mu_2 = E[Y_2], 0 \neq \sigma_1^2 = Var[Y_1], 0 \neq \sigma_2^2 = Var[Y_2]$, and $\sigma_{12} = Cov[Y_1, Y_2]$. Assume that the correlation coefficient $|\rho[Y_1, Y_2]| \neq 1$. Then the joint probability density function of (Y_1, Y_2) , $g : \mathbb{R}^2 \rightarrow [0, \infty)$ is given by

$$g(y_1, y_2) = \frac{\exp\left(-\frac{1}{2(1-\rho^2)} \left[\left(\frac{y_1-\mu_1}{\sigma_1}\right)^2 + \left(\frac{y_2-\mu_2}{\sigma_2}\right)^2 - 2\rho \left(\frac{y_1-\mu_1}{\sigma_1}\right) \left(\frac{y_2-\mu_2}{\sigma_2}\right) \right]\right)}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \quad (6.4.2)$$

Proof- Let a, b be two real numbers. We will show that

$$P(Y_1 \leq a, Y_2 \leq b) = \int_{-\infty}^a \int_{-\infty}^b g(y_1, y_2) dy_2 dy_1. \quad (6.4.3)$$

From the discussion that follows (5.4.1), we can then conclude that the joint density of (Y_1, Y_2) is indeed given by g . To show (6.4.3) we find an alternate description of (Y_1, Y_2) which is the same in distribution. Let Z_1, Z_2 be two independent standard normal random variables. Define

$$\begin{aligned} U &= \sigma_1 Z_1 + \mu_1 \\ V &= \sigma_2(\rho Z_1 + \sqrt{1-\rho^2}Z_2) + \mu_2 \end{aligned} \quad (6.4.4)$$

Let $\alpha, \beta \in \mathbb{R}$. Then

$$\alpha U + \beta V = (\alpha\sigma_1 + \beta\sigma_2\rho)Z_1 + (\beta\sigma_2\sqrt{1-\rho^2})Z_2 + \alpha\mu_1 + \beta\mu_2.$$

As Z_1 and Z_2 are independent standard normal random variables by Theorem 6.3.13, $(\alpha\sigma_1 + \beta\sigma_2\rho)Z_1 + (\beta\sigma_2\sqrt{1-\rho^2})Z_2 \sim \text{Normal}(0, (\alpha\sigma_1 + \beta\sigma_2\rho)^2 + (\beta\sigma_2\sqrt{1-\rho^2})^2)$. Further using Corollary 5.3.3 (a) we have that $\alpha U + \beta V \sim \text{Normal}(\alpha\mu_1 + \beta\mu_2, (\alpha\sigma_1 + \beta\sigma_2\rho)^2 + (\beta\sigma_2\sqrt{1-\rho^2})^2)$. As α, β were arbitrary real numbers by Definition 6.4.1, (U, V) is a bivariate normal random variable.

Using Theorem 6.1.8 and Theorem 6.1.10 (d) that,

$$\mu_1 = E[U], \mu_2 = E[V], \sigma_1^2, \text{ and } \text{Var}[U].$$

Also in addition, using Exercise 6.2.16 and Theorem 6.2.2 (f), we have

$$\begin{aligned} \text{Var}[V] &= \sigma_2^2\rho^2\text{Var}[Z_1] + \sigma_2^2(1-\rho^2)\text{Var}[Z_2] + 2(\sigma_2(\rho + \sqrt{1-\rho^2}))\text{Cov}[Z_1, Z_2] \\ &= \sigma_2^2\rho^2 + \sigma_2^2(1-\rho^2) + 0 = \sigma_2^2 \\ &\text{and} \\ \text{Cov}[U, V] &= \text{Cov}[\sigma_1 Z_1 + \mu_1, \sigma_2(\rho Z_1 + \sqrt{1-\rho^2}Z_2)] \\ &= \sigma_1\sigma_2\rho\text{Cov}[Z_1, Z_1] + \sigma_1\sigma_2\sqrt{1-\rho^2}\text{Cov}[Z_1, Z_2] \\ &= \sigma_1\sigma_2\rho + 0 = \sigma_{12}. \end{aligned}$$

As bivariate normal random variables are by their means and covariances (by Theorem 6.4.2), (Y_1, Y_2) and (U, V) have the same joint distribution. By the above, we have

$$P(Y_1 \leq a, Y_2 \leq b) = P(U \leq a, V \leq b). \quad (6.4.5)$$

By elementary algebra we can also infer from (6.4.4)

$$Z_1 = \frac{U - \mu_1}{\sigma_1}, \quad Z_2 = \frac{V - \mu_2}{\sigma_2\sqrt{1-\rho^2}} - \frac{\rho Z_1}{\sqrt{1-\rho^2}}.$$

So

$$\{U \leq a, V \leq b\} = \left\{ Z_1 \leq \frac{a - \mu_1}{\sigma_1}, Z_2 \leq \frac{b - \mu_2}{\sigma_2 \sqrt{1 - \rho^2}} - \frac{\rho Z_1}{\sqrt{1 - \rho^2}} \right\}$$

So, using this fact in (6.4.5) we get

$$\begin{aligned} P(Y_1 \leq a, Y_2 \leq b) &= P\left(Z_1 \leq \frac{a - \mu_1}{\sigma_1}, Z_2 \leq \frac{b - \mu_2}{\sigma_2 \sqrt{1 - \rho^2}} - \frac{\rho Z_1}{\sqrt{1 - \rho^2}}\right) \\ &= \int_{-\infty}^{\frac{a - \mu_1}{\sigma_1}} \int_{-\infty}^{\frac{b - \mu_2}{\sigma_2 \sqrt{1 - \rho^2}} - \frac{\rho z_1}{\sqrt{1 - \rho^2}}} \frac{\exp(-\frac{z_1^2 + z_2^2}{2})}{2\pi} dz_2 dz_1 \end{aligned} \quad (6.4.6)$$

First performing a u -substitution in the inner integral for each fixed z_1 ,

$$z_2 = \frac{y_2 - \mu_2}{\sigma_2 \sqrt{1 - \rho^2}} - \frac{\rho z_1}{\sqrt{1 - \rho^2}}$$

yields that the inner integral in (6.4.6) for each $z_1 \in \mathbb{R}$

$$\begin{aligned} \int_{-\infty}^{\frac{b - \mu_2}{\sigma_2 \sqrt{1 - \rho^2}} - \frac{\rho z_1}{\sqrt{1 - \rho^2}}} \frac{\exp(-\frac{z_1^2 + z_2^2}{2})}{2\pi} dz_2 &= \int_{-\infty}^b \frac{\exp(-\frac{z_1^2 + (\frac{y_2 - \mu_2}{\sigma_2 \sqrt{1 - \rho^2}} - \frac{\rho z_1}{\sqrt{1 - \rho^2}})^2}{2})}{2\pi \sigma_2 \sqrt{1 - \rho^2}} dy_2 \\ &= \int_{-\infty}^b \frac{\exp(-\frac{1}{2(1 - \rho^2)}[(1 - \rho^2)z_1^2 + (\frac{y_2 - \mu_2}{\sigma_2} - \rho z_1)^2])}{2\pi \sigma_2 \sqrt{1 - \rho^2}} dy_2 \\ &= \int_{-\infty}^b \frac{\exp(-\frac{1}{2(1 - \rho^2)}[z_1^2 + (\frac{y_2 - \mu_2}{\sigma_2})^2 - 2\rho(\frac{y_2 - \mu_2}{\sigma_2})z_1])}{2\pi \sigma_2 \sqrt{1 - \rho^2}} dy_2. \end{aligned}$$

Substituting the above into (6.4.6), we have

$$P(Y_1 \leq a, Y_2 \leq b) = \int_{-\infty}^{\frac{a - \mu_1}{\sigma_1}} \int_{-\infty}^b \frac{\exp(-\frac{1}{2(1 - \rho^2)}[z_1^2 + (\frac{y_2 - \mu_2}{\sigma_2})^2 - 2\rho(\frac{y_2 - \mu_2}{\sigma_2})z_1])}{2\pi \sigma_2 \sqrt{1 - \rho^2}} dy_2 dz_1 \quad (6.4.7)$$

Performing a u -substitution

$$z_1 = \frac{y_1 - \mu_1}{\sigma_1}$$

on the outer integral above we obtain

$$\begin{aligned} P(Y_1 \leq a, Y_2 \leq b) &= \int_{-\infty}^a \int_{-\infty}^b \frac{\exp(-\frac{1}{2(1 - \rho^2)}[(\frac{y_1 - \mu_1}{\sigma_1})^2 + (\frac{y_2 - \mu_2}{\sigma_2})^2 - 2\rho(\frac{y_2 - \mu_2}{\sigma_2})(\frac{y_1 - \mu_1}{\sigma_1})])}{2\pi \sigma_1 \sigma_2 \sqrt{1 - \rho^2}} dy_2 dy_1 \end{aligned}$$

Thus we have established (6.4.3). ■

EXERCISES

Ex. 6.4.1. Let X_1, X_2 be two independent Normal random variables with mean 0 and variance 1. Show that (X_1, X_2) is a bivariate normal random variable.

Ex. 6.4.2. Let (X_1, X_2) be a bivariate normal random variable. Assume that the correlation coefficient $|\rho[X_1, X_2]| \neq 1$. Show that X_1 and X_2 are Normal random variables by calculating their marginal densities.

Ex. 6.4.3. Let X_1, X_2 be two independent normal random variables with mean 0 and variance 1. Let (Y_1, Y_2) be a bivariate normal random variable with zero means, variances equal to 1 and correlation $\rho = \rho[Y_1, Y_2]$, with $\rho^2 \neq 1$. Let f be the joint probability density function of (X_1, X_2) and g be the joint probability density function of (Y_1, Y_2) . For $0 < \alpha < 1$, let (Z_1, Z_2) be a bivariate random variable with joint density given by

$$h(z_1, z_2) = \alpha g(z_1, z_2) + (1 - \alpha)f(z_1, z_2),$$

for any real numbers z_1, z_2 .

- (a) Write down the exact expressions for f and g .
- (b) Verify that h is indeed a probability density function.
- (c) Show that Z_1 and Z_2 are Normal random variables by calculating their marginal densities.
- (d) Show that (Z_1, Z_2) is not a bivariate normal random variable.

Ex. 6.4.4. Suppose X_1, X_2, \dots, X_n are independent and normally distributed. Let $Y = c_1 X_1 + \dots + c_n X_n$ and let $Z = d_1 X_1 + \dots + d_n X_n$ be linear combinations of these variables (for real numbers c_j and d_j). Then (Y, Z) is bivariate normal.

Ex. 6.4.5. Prove Theorem 6.3.13. Specifically, suppose for $i = 1, 2, \dots, n$ that $X_i \sim \text{Normal}(\mu_i, \sigma_i^2)$ with X_1, X_2, \dots, X_n independent. Let a_1, a_2, \dots, a_n be real numbers, not all zero, and let $Y = a_1 X_1 + a_2 X_2 + \dots + a_n X_n$. Prove that Y is normally distributed and find its mean and variance in terms of the a 's, μ 's, and σ 's.

Ex. 6.4.6. Let (X_1, X_2) be a bivariate Normal random variable. Define

$$\Sigma = \begin{bmatrix} \text{Cov}[X_1, X_1] & \text{Cov}[X_1, X_2] \\ \text{Cov}[X_2, X_1] & \text{Cov}[X_2, X_2] \end{bmatrix}$$

$$\text{and } \mu_1 = E[X_1], \mu_2 = E[X_2], \mu_{2 \times 1} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}.$$

Σ is referred to as the covariance matrix of (X_1, X_2) and μ is the mean matrix of (X_1, X_2) .

- (a) Compute $\det(\Sigma)$.
- (b) Show that the joint density of (X_1, X_2) can be rewritten in matrix notation as

$$g(x_1, x_2) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_2 \end{bmatrix} \Sigma^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}\right)$$

(c)

$$A_{2 \times 2} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \eta_{2 \times 1} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}$$

such that a_{ij} are real numbers. Suppose we define

$$Y = AX = \begin{bmatrix} a_{11}X_1 + a_{12}X_2 + \eta_1 \\ a_{21}X_1 + a_{22}X_2 + \eta_2 \end{bmatrix}.$$

Then (Y_1, Y_2) is also a bivariate Normal random variable, with covariance matrix $A\Sigma A^T$ and mean matrix $A\mu + \eta$.

Hint: Compute means, variances and covariances of Y_1, Y_2 and use Theorem 6.4.2

7

SAMPLING AND DESCRIPTIVE STATISTICS

The distinction between Probability and Statistics is somewhat fuzzy, but largely has to do with the perspective of what is known versus what is to be determined. One may think of Probability as the study of models for (random) experiments when the model is fully known. When the model is not fully known and one tries to infer about the unknown aspects of the model based on observed outcomes of the experiment, this is where Statistics enters the picture. In this chapter we will be interested in problems where we assume we know the outputs of random variables, and wish to use that information to say what we can about their (unknown) distributions.

Suppose, for instance, we sample from a large population and record a numerical fact associated with each selection. This may be recording the heights of people, recording the arsenic content of water samples, recording the diameters of randomly selected trees, or anything else that may be thought of as repeated, random measurements. Sampling an individual from a population in this case may be viewed as a random experiment. If the sampling were done at random with replacement with each selection independent of any other, we could view the resulting numerical measurements as i.i.d. random variables X_1, X_2, \dots, X_n . A more common situation is sampling without replacement, but we have previously seen (See Section 2.3) that when the sample size is small relative to the size of the population, the two sampling methods are not dramatically different. In this case we have the results of n samples from a distribution, but we don't actually know the distribution itself. How might we use the samples to attempt to predict such things as expected value and variance?

7.1 THE EMPIRICAL DISTRIBUTION

A natural quantity we can create from the observed data, regardless of the underlying distribution that generated it, is a discrete distribution that puts equal probability on each observed point. This distribution is known as the empirical distribution. Some values of X_i can of course be repeated, so the empirical distribution is formally defined as follows.

DEFINITION 7.1.1. Let X_1, X_2, \dots, X_n be i.i.d. random variables. The “empirical distribution” based on these is the discrete distribution with probability mass function given by

$$f(t) = \frac{1}{n} \#\{X_i = t\}.$$

We can now study the empirical distribution using the tools of probability. Doing so does not make any additional assumptions about the underlying distribution, and inferences about it based on the empirical distribution are traditionally referred to as “descriptive statistics”. In later chapters, we will see that making additional assumptions lets us make “better” inferences, provided the additional assumptions are valid.

It is important to realize that the empirical distribution is itself a random quantity, as each sample realisation will produce a different discrete distribution. We intuitively expect it to carry information about the underlying distribution, especially as the sample size n grows. For example, the expectation computed from the empirical distribution should be closely related to the true underlying expectation, probabilities of events computed from the empirical distribution should be related to the true probabilities of those events, and so on. In the remainder of this chapter, we will make this intuition more precise and describe some tools to investigate the properties of the empirical distribution.

7.2 DESCRIPTIVE STATISTICS

7.2.1 Sample Mean

Given a sample of observations, we define the sample mean to be the familiar definition of average.

DEFINITION 7.2.1. Let X_1, X_2, \dots, X_n be i.i.d. random variables. The “sample mean” of these is

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}.$$

It is easy to see that \bar{X} is the expected value of a random variable whose distribution is the empirical distribution based on X_1, X_2, \dots, X_n (see Exercise 7.2.4). Suppose the X_j random variables have a finite expected value μ . The sample mean \bar{X} is not the same as this expected value. In particular μ is a fixed constant while \bar{X} is a random variable. From the statistical perspective, μ is usually assumed to be an unknown quantity while \bar{X} is something that may be computed from the results of the sample X_1, X_2, \dots, X_n . How well does \bar{X} work as an estimate of μ ? The next theorem begins to answer this question.

THEOREM 7.2.2. Let X_1, X_2, \dots, X_n be an i.i.d. sample of random variables whose distribution has finite expected value μ and finite variance σ^2 . Let \bar{X} represent the sample mean. Then

$$E[\bar{X}] = \mu \quad \text{and} \quad SD[\bar{X}] = \frac{\sigma}{\sqrt{n}}.$$

Proof -

$$\begin{aligned} E[\bar{X}] &= E\left[\frac{X_1 + X_2 + \dots + X_n}{n}\right] \\ &= \frac{E[X_1] + E[X_2] + \dots + E[X_n]}{n} \\ &= \frac{n\mu}{n} = \mu \end{aligned}$$

To calculate the standard deviation, we consider the variance and use Theorem 4.2.6 and Exercise 6.1.12 to obtain

$$\begin{aligned} Var[\bar{X}] &= Var\left[\frac{X_1 + X_2 + \dots + X_n}{n}\right] \\ &= \frac{Var[X_1] + Var[X_2] + \dots + Var[X_n]}{n^2} \\ &= \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \end{aligned}$$

Taking square roots then shows $SD[\bar{X}] = \frac{\sigma}{\sqrt{n}}$. ■

The fact that $E[\bar{X}] = \mu$ means that, on average, the quantity \bar{X} is accurately describing the unknown mean μ . In the language of statistics \bar{X} is said to be an “unbiased estimator” of the quantity μ . Note also that $SD[\bar{X}] \rightarrow 0$ as $n \rightarrow \infty$ meaning that the larger the sample size, the more accurately \bar{X} reflects its average of μ . In other words, if there is an unknown distribution from which it is possible to sample, averaging a large sample should produce a value close to the expected value of the distribution. In technical terms, this means that the sample mean is a “consistent estimator” of the population mean μ .

7.2.2 Sample Variance

Given a sample of observations, we define the sample variance below.

DEFINITION 7.2.3. Let X_1, X_2, \dots, X_n be i.i.d. random variables. The “sample variance” of these is

$$S^2 = \frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n - 1}.$$

Note that this definition is not universal; it is common to define sample variance with n (instead of $n - 1$) in the denominator, in which case the definition matches the variance of the empirical distribution of X_1, X_2, \dots, X_n (Exercise 7.2.4). The definition given here produces a quantity that is unbiased for the underlying population variance, a fact that follows from the next theorem.

THEOREM 7.2.4. Let X_1, X_2, \dots, X_n be an i.i.d. sample of random variables whose distribution has finite expected value μ and finite variance σ^2 . Then S^2 is an unbiased estimator of σ^2 , i.e.

$$E[S^2] = \sigma^2.$$

Proof - First note that

$$E[\bar{X}^2] = \text{Var}[\bar{X}] + (E[\bar{X}])^2 = \frac{\sigma^2}{n} + \mu^2$$

whereas

$$E[X_j^2] = \text{Var}[X_j] + E[X_j]^2 = \sigma^2 + \mu^2.$$

Now consider the quantity $(n - 1)S^2$.

$$\begin{aligned} E[(n - 1)S^2] &= E[(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2] \\ &= E[X_1^2 + X_2^2 + \dots + X_n^2] - 2E[(X_1 + X_2 + \dots + X_n)\bar{X}] \\ &\quad + E[\bar{X}^2 + \bar{X}^2 + \dots + \bar{X}^2] \end{aligned}$$

But $X_1 + X_2 + \dots + X_n = n\bar{X}$, so

$$\begin{aligned} E[(n - 1)S^2] &= E[X_1^2 + X_2^2 + \dots + X_n^2] - 2nE[\bar{X}^2] + nE[\bar{X}^2] \\ &= E[X_1^2 + X_2^2 + \dots + X_n^2] - nE[\bar{X}^2] \\ &= n(\sigma^2 + \mu^2) - n\left(\frac{\sigma^2}{n} + \mu^2\right) \\ &= (n - 1)\sigma^2 \end{aligned}$$

Dividing by $n - 1$ gives the desired result, $E[S^2] = \sigma^2$. ■

A more important property (than unbiasedness) is that S^2 and its variant with n in the denominator are both “consistent” for σ^2 , just as \bar{X} was for μ , in the sense that $\text{Var}[S^2] \rightarrow 0$ as $n \rightarrow \infty$ under some mild conditions.

7.2.3 Sample proportion

Expectation and variance are commonly used summaries of a random variable, but they do not characterize its distribution completely. In general, the distribution of a random variable X is fully known if we can compute $P(X \in A)$ for any event A . In particular, it is enough to know probabilities of the type $P(X \leq t)$, which is precisely the cumulative distribution function of X evaluated at t .

Given a sample of i.i.d. observations X_1, X_2, \dots, X_n from a common distribution defined by a random variable X , the probability $P(X \in A)$ of any event A has the natural sample analog $P(Y \in A)$, where Y

is a random variable following the empirical distribution based on X_1, X_2, \dots, X_n . To understand this quantity, recall that Y essentially takes values X_1, X_2, \dots, X_n with probability $1/n$ each, and so we have

$$P(Y \in A) = \sum_{X_i \in A} \frac{1}{n} = \frac{\#\{X_i \in A\}}{n}$$

In other words, $P(Y \in A)$ is simply the proportion of sample observations for which the event A happened. Not surprisingly, $P(Y \in A)$ is a good estimator of $P(X \in A)$ in the following sense.

THEOREM 7.2.5. *Let X_1, X_2, \dots, X_n be an i.i.d. sample of random variables with the same distribution as a random variable X , and suppose that we are interested in the value $p = P(X \in A)$ for an event A . Let*

$$\hat{p} = \frac{\#\{X_i \in A\}}{n}.$$

Then, $E(\hat{p}) = P(X \in A)$ and $\text{Var}(\hat{p}) \rightarrow 0$ as $n \rightarrow \infty$.

Proof - Let

$$Y = \#\{X_i \in A\} = \sum_{i=1}^n Z_i,$$

where for $1 \leq i \leq n$,

$$Z_i = \begin{cases} 1 & \text{if } X_i \in A \\ 0 & \text{otherwise} \end{cases}$$

It is easy to see $P(Z_i = 1) = P(X_i \in A) = p$. Further Z_i 's are independent because X_i 's are independent (See Theorem 3.3.6 and Exercise 7.2.1). Thus, Y has the Binomial distribution with parameters n and p , with expectation np and variance $np(1-p)$. It immediately follows that

$$E(\hat{p}) = E(Y/n) = p \text{ and } \text{Var}(\hat{p}) = p(1-p)/n$$

which has the limiting value 0 as $n \rightarrow \infty$. ■

This result is a special case of the more general “law of large numbers” we will encounter in Section 8.2. It is important because it gives formal credence to our intuition that the probability of an event measures the limiting relative frequency of that event over repeated trials of an experiment.

DEFINITION 7.2.6. *In terms of our notation above, the analog of the cumulative distribution function of X is the cumulative distribution function of Y , which is traditionally denoted by*

$$\hat{F}_n(t) = P(Y \leq t) = \frac{\#\{X_i \leq t\}}{n}$$

and known as the “empirical cumulative distribution function” or ECDF of X_1, X_2, \dots, X_n .

EXERCISES

Ex. 7.2.1. Verify that the proofs of Theorem 3.3.5 and Theorem 3.3.6 hold for continuous random variables.

Ex. 7.2.2. Let X and Y be two continuous random variables having the same distribution. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a piecewise continuous function. Then show that $f(X)$ and $f(Y)$ have the same distribution.

Ex. 7.2.3. Verify Exercise 7.2.2 for discrete random variables.

Ex. 7.2.4. Let P be the empirical distribution defined by sample observations X_1, X_2, \dots, X_n . In other words, P is the discrete distribution with probability mass function given in Definition 7.1.1. Let Y be a random variable with distribution P .

- (a) Show that $E(Y) = \bar{X}$.
- (b) Show that $Var(Y) = \frac{n-1}{n}S^2$.

Ex. 7.2.5. Let X_1, X_2, \dots, X_n be i.i.d. random variables with finite expectation μ , finite variance σ^2 , and finite $\gamma = E(X_1 - \mu)^4$. Compute $Var(S^2)$ in terms of μ , σ^2 , and γ and show that $Var(S^2) \rightarrow 0$ as $n \rightarrow \infty$.

Ex. 7.2.6. Let X_1, X_2, \dots, X_n be i.i.d. random variables with finite expectation μ and finite variance σ^2 . let $S = \sqrt{S^2}$, the non-negative root of the sample variance. The quantity S is called the “sample standard deviation”. Although $E[S^2] = \sigma^2$, it is not true that $E[S] = \sigma$. In other words, S is not an unbiased estimator for σ . Follow the steps below to see why.

- (a) Let Z be a random variable with finite mean and finite variance. Prove that $E[Z^2] \geq E[Z]^2$ and give an example to show that equality may not hold. (Hint: Consider how these quantities relate to the variance of Z).
- (b) Use (a) to explain why $E[S] \leq \sigma$ and give an example to show that equality may not hold.

7.3 SIMULATION

The preceding discussion gives several mathematical statements about random samples, but it is difficult to develop any intuition about what these statements mean unless we look at actual data. Data is of course abundant in our world, and we will look at some real life data sets later in this book. However, the problem with real data is that we do not usually know for certain the random variable that generated it. To hone our intuition, it is therefore useful to be able to generate random samples from a distribution we specify. The process of doing so using a computer program is known as “simulation”.

Simulation is not an easy task, because computers are by nature not random. Simulation is in fact not a random process at all; it is a completely deterministic process that tries to mimic randomness. We will not go into how simulation is done, but simply use R to obtain simulated random samples.

R supports simulation from many distributions, including all the ones we have encountered. The general pattern of usage is that each distribution has a corresponding function that is called with the sample size as an argument, and further arguments specifying parameters. The function returns the simulated observations as a vector. For example, 30 Binomial(100, 0.75) samples can be generated by

```
> rbinom(30, size = 100, prob = 0.75)
[1] 74 84 87 75 69 71 80 75 79 68 72 75 78 75 76 78 82 70 74 76 74 77 70 73 76
[26] 70 70 76 72 77
```

We usually want to do more than just print simulated data, so we typically store the result in a variable and make further calculations with it; for example, compute the sample mean, or the sample proportion of cases where a particular event happens.

```
> x <- rbinom(30, size = 100, prob = 0.75)
> mean(x)
```

```
[1] 73.63333
> sum(x >= 75) / length(x)
[1] 0.4333333
```

R has a useful function called `replicate` that allows us to repeat such an experiment several times.

```
> replicate(15, {
+   x <- rbinom(30, size = 100, prob = 0.75)
+   mean(x)
+ })
```

```
[1] 73.23333 75.53333 74.50000 75.46667 75.36667 75.63333 73.53333 74.66667
[9] 75.43333 73.96667 74.40000 75.16667 74.40000 74.16667 75.30000
```

```
> replicate(15, {
+   x <- rbinom(30, size = 100, prob = 0.75)
+   sum(x >= 75) / length(x)
+ })
[1] 0.5000000 0.4333333 0.8666667 0.5333333 0.6000000 0.5000000 0.5666667
[8] 0.5333333 0.6333333 0.4000000 0.5333333 0.5333333 0.5666667 0.5333333
[15] 0.4666667
```

This gives us an idea of the variability of the sample mean and sample proportion computed from a sample of size 30. We know of course that the sample mean has expectation $100 \times 0.75 = 75$, and we can compute the expected value of the proportion using R as follows.

```
> 1 - pbinom(74, size = 100, prob = 0.75)
[1] 0.5534708
```

So the corresponding estimates are close to the expected values, but with some variability. We expect the variability to go down if the sample size increases, say, from 30 to 3000.

```
> replicate(15, {
+   x <- rbinom(3000, size = 100, prob = 0.75)
+   mean(x)
+ })
[1] 75.00300 75.11233 74.95167 74.99033 75.06167 74.96633 74.86000 74.94633
[9] 75.08333 74.92700 75.03167 75.02633 75.05000 74.95467 75.03167

> replicate(15, {
+   x <- rbinom(3000, size = 100, prob = 0.75)
+   sum(x >= 75) / length(x)
+ })
[1] 0.5706667 0.5780000 0.5433333 0.5440000 0.5863333 0.5426667 0.5496667
[8] 0.5440000 0.5516667 0.5486667 0.5423333 0.5480000 0.5526667 0.5403333
[15] 0.5573333
```

Indeed we see that the estimates are much closer to their expected values now.

We can of course replicate this process for other events of interest, and indeed for many other distributions. We will see in the next section how we can simulate observations following the normal distribution using the function `rnorm`, and the exponential distribution using the function `rexp`. It is also interesting to think about how one can simulate observations from a given distribution when a function to do so is not already available. The following examples explore some simple approaches.

EXAMPLE 7.3.1. When trying to formulate a method to simulate random variables from a new distribution, it is customary to assume that we already have a method to generate random variables from $\text{Uniform}(0, 1)$. Let us see this can be used to generate random observations from a $\text{Poisson}(\lambda)$ distribution using its probability mass function.

Let X denote an observation from the $\text{Poisson}(\lambda)$ distribution, and $U \sim \text{Uniform}(0, 1)$. Denote $p_i = P(X = i)$. The basic idea is as follows:

$$\begin{aligned} p_0 &= P(U \leq p_0) \\ P(U \leq p_0 + p_1) &= p_0 + p_1 \Rightarrow p_1 = P(p_0 < U < p_0 + p_1) \\ P(U \leq p_0 + p_1 + p_2) &= p_0 + p_1 + p_2 \Rightarrow p_2 = P(p_0 + p_1 < U < p_0 + p_1 + p_2) \end{aligned}$$

and so on. Thus, if we set Y to be 0 if $U \leq p_0$, and k if U satisfies $\sum_{i=0}^{k-1} p_i < U < \sum_{i=0}^k p_i$, then Y has the same distribution as X .

To use this idea to generate 50 observations from Poisson(5), we can use the following code in R, noting that $\sum_{i=0}^k p_i = P(X \leq k)$.

```
> replicate(50,
+           {
+             U <- runif(1)
+             Y <- 0
+             while (U > ppois(Y, lambda = 5)) Y <- Y + 1
+             Y
+           })
[1] 4 8 3 4 7 3 7 5 3 4 5 1 2 8 6 3 4 2 8 5 7 2 4 4
[26] 3 4 8 5 6 8 3 7 9 5 5 5 7 8 4 5 3 3 8 2 8 2 7 8 14
```

Of course, there is nothing in this procedure that is specific to the Poisson distribution. By replacing the call to `ppois()` suitably, the same process can be used to simulate random observations from any discrete distribution supported on the non-negative integers. ■

EXAMPLE 7.3.2. The process described in the previous example cannot be used for continuous random variables. In such cases, Lemma 5.3.7 often proves useful. The first part of the lemma states that if $U \sim \text{Uniform}(0, 1)$, and F_X is the distribution function of a continuous random variable X , then $Y = F_X^{-1}(U)$ has the same distribution as X . This can be used to generate observations from X provided we can compute F_X^{-1} .

Consider the case where we want X to have the $\text{Exp}(1)$ distribution. Then, $F_X(x) = 1 - e^{-x}$ for $x > 0$. Solving for $F_X(x) = u$, we have

$$\begin{aligned} 1 - e^{-x} &= u \\ \Rightarrow e^{-x} &= 1 - u \\ \Rightarrow x &= -\log(1 - u), \end{aligned}$$

that is, $F_X^{-1}(u) = -\log(1 - u)$. Thus, we can simulate 50 observations from the $\text{Exp}(1)$ distribution using the following R code.

```
> -log(1 - runif(50))
[1] 0.17983033 0.59899225 0.39765691 0.46661641 1.83186881 0.75753630
[7] 0.15224550 3.01323320 0.02324019 2.62589324 0.50319325 0.06495110
[13] 1.73626921 0.79253356 0.46701605 1.31246443 1.94788764 0.32681347
[19] 0.96975851 0.52949759 0.74217408 0.85115821 0.04679527 0.35540345
[25] 0.25261271 0.91725848 0.54630522 1.53183895 0.52956653 1.02305166
[31] 1.65161608 1.30340256 0.27096431 1.05641695 0.58749136 0.19851994
[37] 0.04194768 1.43645222 0.70200050 1.09493028 0.40181847 1.76807864
[43] 3.24628447 0.65443582 0.08138553 1.23594540 0.28568794 1.90748439
[49] 0.27814493 0.54204644
```

multiple values at once, and the fact that the expression for $F_X^{-1}(u)$ can be easily vectorized. We can multiply the resulting observations by $1/\lambda$ to simulate observations from the $\text{Exp}(\lambda)$ distribution. ■

EXAMPLE 7.3.3. (TODO: Simulate Bivariate Normal) Suppose we want to simulate observations (X, Y) from a bivariate normal distribution. To start with, let us assume that both mean parameters are 0, both variance parameters are 1, and the correlation coefficient ρ (which is also the covariance) is specified.

This problem is somewhat tricky, because the definition of bivariate normal does not directly provide a way to simulate it. All we know is that any linear combination $aX + bY$ has a univariate normal distribution. ■

EXERCISES

Ex. 7.3.1. (a) Show that both the sample mean and the sample variance of a sample obtained from the $\text{Poisson}(\lambda)$ distribution will be unbiased estimators of λ .

- (b) Which of these estimators is better? To answer this question, simulate random observations from the $\text{Poisson}(\lambda)$ distribution for various values of λ using the R function `rpois`. Explore the behaviour of the two estimates by varying λ as well as the sample size.

Ex. 7.3.2. Exercise 2.3.7 described the technique called “capture-recapture” which biologists use to estimate the size of the population of a species that cannot be directly counted. Suppose the unknown population size is N , and fifty members of the species are selected and given an identifying mark. Sometime later a sample of size twenty is taken from the population, and it is found to contain X of the twenty previously marked. Equating the proportion of marked members in the second sample and the population, we have $\frac{X}{20} = \frac{50}{N}$, giving an estimate of $\hat{N} = \frac{1000}{X}$.

Recall that X has a hypergeometric distribution that involves N as a parameter. It is not easy to compute $E[\hat{N}]$ and $\text{Var}[\hat{N}]$; however, Hypergeometric random variables can be simulated in R using the function `rhyper`. For each $N = 50, 100, 200, 300, 400$, and 500 , use this function to simulate 1000 values of \hat{N} and use them to estimate $E[\hat{N}]$ and $\text{Var}[\hat{N}]$. Plot these estimates as a function of N .

Ex. 7.3.3. Suppose p is the unknown probability of an event A , and we estimate p by the sample proportion \hat{p} based on an i.i.d. sample of size n .

- (a) Write $\text{Var}[\hat{p}]$ and $SD[\hat{p}]$ as functions of n and p .
- (b) Using the relations derived above, determine the sample size n , as a function of p , that is required to achieve $SD(\hat{p}) = 0.01$. How does this required value of n vary with p ?
- (c) Design and implement the following simulation study to verify this behaviour. For $p = 0.01, 0.1, 0.25, 0.5, 0.75, 0.9$, and 0.99 ,
 - (i) Simulate 1000 values of \hat{p} with $n = 500$.
 - (ii) Simulate 1000 values of \hat{p} with n chosen according to the formula derived above.

In each case, you can think of the 1000 values as i.i.d. samples from the distribution of \hat{p} , and use the sample standard deviation as an estimate of $SD[\hat{p}]$. Plot the estimated values of $SD(\hat{p})$ against p for both choices of n . Your plot should look similar to Figure 7.1.

- (d) (FIXME: Open-ended question) Do you think the standard deviation $SD[\hat{p}]$ is a good way to measure how well \hat{p} measures p ? If not, what alternatives can you think of?

Ex. 7.3.4. TODO: Give several other distributions as specific examples and specific events. Mention corresponding R functions.

7.4 PLOTS

As we will see in later chapters, making more assumptions about the underlying distribution of X allows us to give concrete answers to many important questions. This is indeed a standard and effective approach to doing statistics, but in following that approach there is a danger of forgetting that assumptions have been made, which we should guard against by doing our best to convince ourselves beforehand that the assumptions we are making are reasonable.

Doing this is more of an art than a science, and usually takes the form of staring at plots obtained from the sample observations, with the hope of answering the question: “does this plot look like what I would have expected it to look like had my assumptions been valid?” Remember that the sample X_1, X_2, \dots, X_n is a random sample, so any plot derived from it is also a “random plot”. Unlike simple quantities such as sample mean and sample variance, it is not clear what to “expect” such plots to look like, and the only way to really hone our instincts to spot anomalies is through experience. In this section, we introduce some commonly used plots and use simulated data to give examples of how such plots might look like when the usual assumptions we make are valid or invalid.

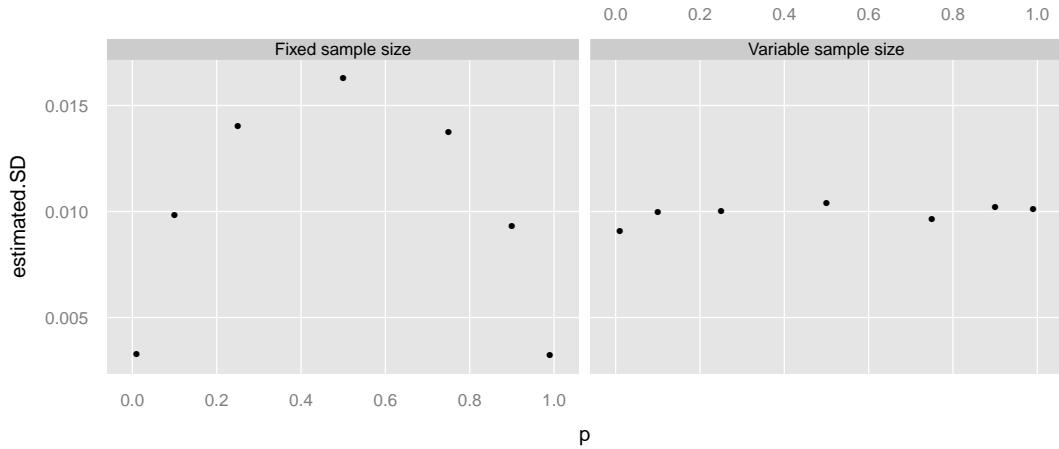


Figure 7.1: Estimated standard deviation in estimating a probability using sample proportion as a function of the probability being estimated. See exercise 7.3.3.

7.4.1 Empirical Distribution Plot for Discrete Distributions

The typical assumption made about a random sample is that the underlying random variable belongs to a *family* of distributions rather than a very specific one. For example, we may assume that the random variable has a Poisson(λ) distribution for some $\lambda > 0$, without placing any further restriction on λ , or a Binomial(n, p) distribution for some $0 < p < 1$. Such families are known as parametric families.

When the data X_1, X_2, \dots, X_n are from a discrete distribution, the simplest representation of the data is its empirical distribution, which is essentially a table of the frequencies of each value that appeared. For example, if we simulate 1000 samples from a Poisson distribution with mean 3, its frequency table may look like

```
> x <- rpois(1000, lambda = 3)
> table(x)

x
 0   1   2   3   4   5   6   7   8   9 
 56 154 206 236 153 111  48  21  13  2 

> prop.table(table(x))

x
 0   1   2   3   4   5   6   7   8   9 
0.056 0.154 0.206 0.236 0.153 0.111 0.048 0.021 0.013 0.002
```

The simplest graphical representation of such a table is through a plot similar to Figure 7.2, which represents a larger Poisson sample with mean 30, resulting in many more distinct values. Although in theory all non-negative integers have positive probability of occurring, the probabilities are too small to be relevant beyond a certain range. This plot does not have a standard name, although it may be considered a variant of the Cleveland Dot Plot. We will refer to it as the *Empirical Distribution Plot* from now on.

We can make similar plots for samples from Binomial or any other distribution. Unfortunately, looking at this plot does not necessarily tell us whether the underlying distribution is Poisson, in part because the shape of the Poisson distribution varies with the λ parameter. A little later, We will discuss a modification of the empirical distribution plot, known as a rootogram, that helps make this kind of comparison a little easier.

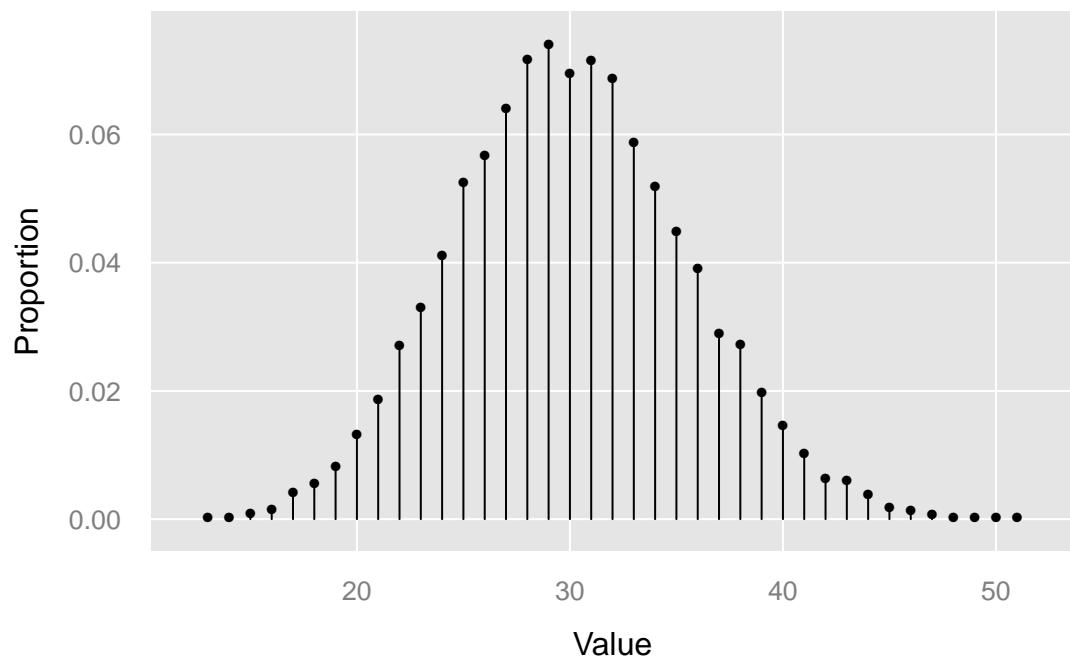


Figure 7.2: Empirical frequency distribution of 10000 random samples from the $\text{Poisson}(30)$ distribution.

7.4.2 Histograms for Continuous Distributions

In the case of continuous distributions, we similarly want to make assumptions about a random sample being from a parametric family of distributions. For example, we may assume that the random variable has a $\text{Normal}(\mu, \sigma^2)$ distribution without placing any further restriction on the parameters μ or σ^2 (except of course that $\sigma^2 > 0$), or that it has an $\text{Exponential}(\lambda)$ distribution with any value of the parameter $\lambda > 0$. Such families are known as parametric families. For both these examples, the *shape* of the distribution does not depend on the parameters, and this makes various diagnostic plots more useful.

The empirical distribution plot above is not useful for data from a continuous distribution, because by the very nature of continuous distributions, all the data points will be distinct with probability 1, and the value of the empirical distribution function will be exactly $1/n$ at these points.

The plot that is most commonly used instead to study distributions is the histogram. It is similar to the empirical distribution plot, except that it does not retain all the information contained in the empirical distribution, and instead divides the range of the data into arbitrary bins and counts the frequencies of data points falling into each bin. More precisely, the histogram estimates the probability density function of the underlying random variable by estimating the density in each bin as a quantity such that the probability of each bin is proportional to the number of observations in that bin. By choosing the bins judiciously, for example by having more of them as sample size increases, the histogram strikes a balance that ensures that the histogram “converges” to the true underlying density as $n \rightarrow \infty$.

Figure 7.3 gives examples of histograms where data are simulated from the normal and exponential distributions for varying sample sizes. Five replications are shown for each sample size. We can see that for large sample sizes, the shape of the histograms are recognizably similar to the shapes of the corresponding theoretical distributions seen in Figure 5.1 and Figure 5.2 in Chapter 5. Moreover, the shape is consistent over the five replications. This is not true, however, for small sample sizes. Remember that the histograms are based on the observed data, and are therefore random objects themselves. As we saw with numerical properties like the mean, estimates have higher variability when the sample size is small, and get less variable as sample size increases. The same holds for graphical estimates, although making this statement precise is more difficult.

7.4.3 Hanging Rootograms for Comparing with Theoretical Distributions

Graphical displays of data are almost always used for some kind of comparison. Sometimes these are implicit comparisons, say, asking how many peaks does a density have, or is it symmetric? More often, they are used to compare samples from two subpopulations, say, the distribution of height in males and females. Sometimes, as discussed above, they are used to compare an observed sample to a hypothesized distribution.

In the case of the empirical distribution plot, a simple modification is to add the probability mass function of the theoretical distribution. This, although a reasonable modification, is not optimal. Research into human perception of graphical displays indicates that the human eye is more adept at detecting departures from straight lines than from curves. Taking this insight into account, John Tukey suggested “hanging” the vertical lines in an empirical distribution plot (which are after all nothing but sample proportions) from their expected values under the hypothesized distribution. He further suggested a transformation of what is plotted: instead of the sample proportions and the corresponding expected probabilities, he suggested plotting their square roots, thus leading to the name *hanging rootogram* for the resulting plot. The reason for making this transformation is as follows. Recall that for a proportion \hat{p} obtained from a sample of size n ,

$$\text{Var}[\hat{p}] = \frac{p(1-p)}{n} \approx \frac{p}{n}$$

provided p is close to 0. In Chapter 9, we will encounter the Central Limit Theorem and the Delta Method, which can be used to show that as the sample size n grows large, $\text{Var}[\sqrt{\hat{p}}] \approx c/n$ for a constant c . This means that unlike $\hat{p} - p$, the variance of $\sqrt{\hat{p}} - \sqrt{p}$ will be approximately independent of p . Figure 7.4 gives examples of hanging rootograms.

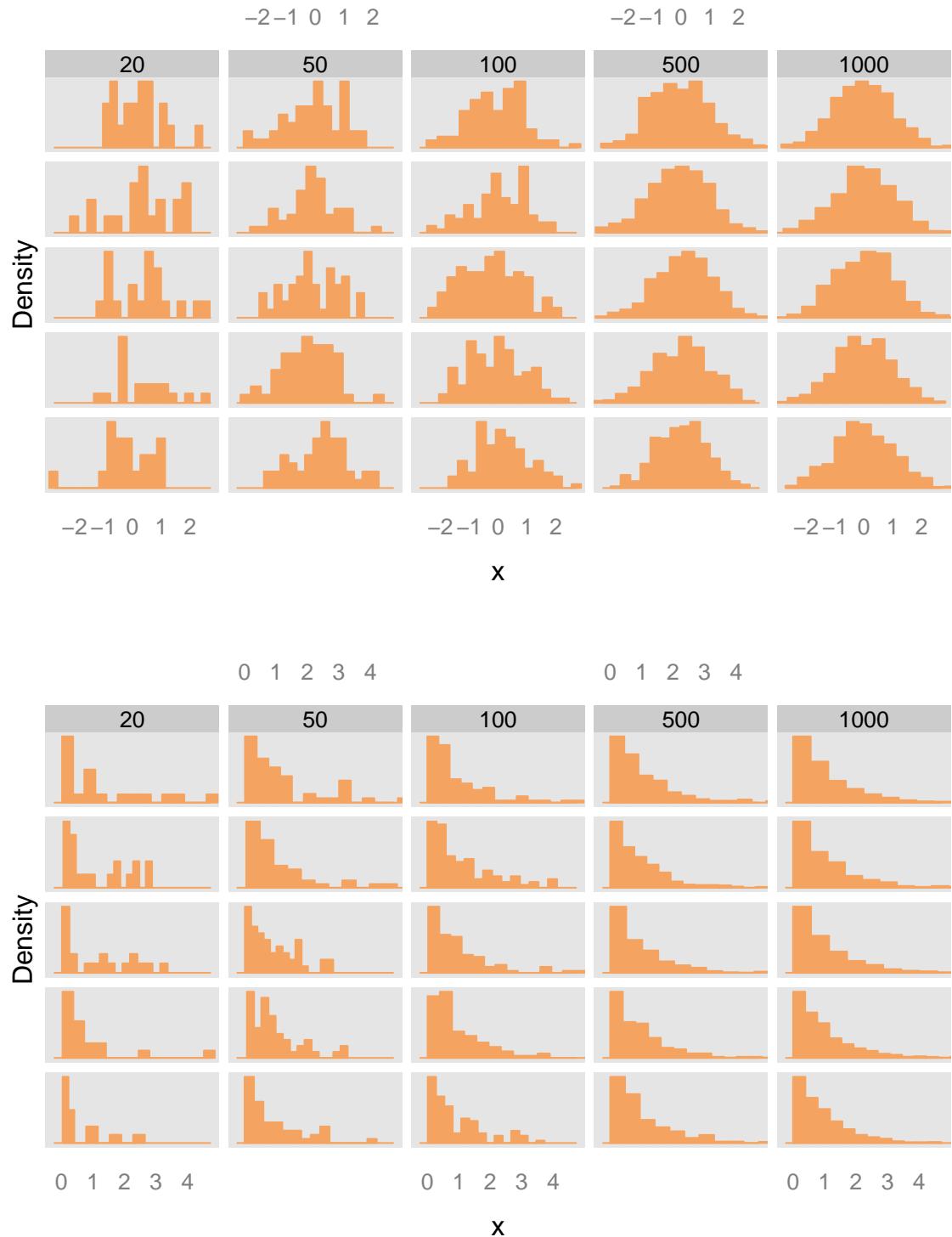
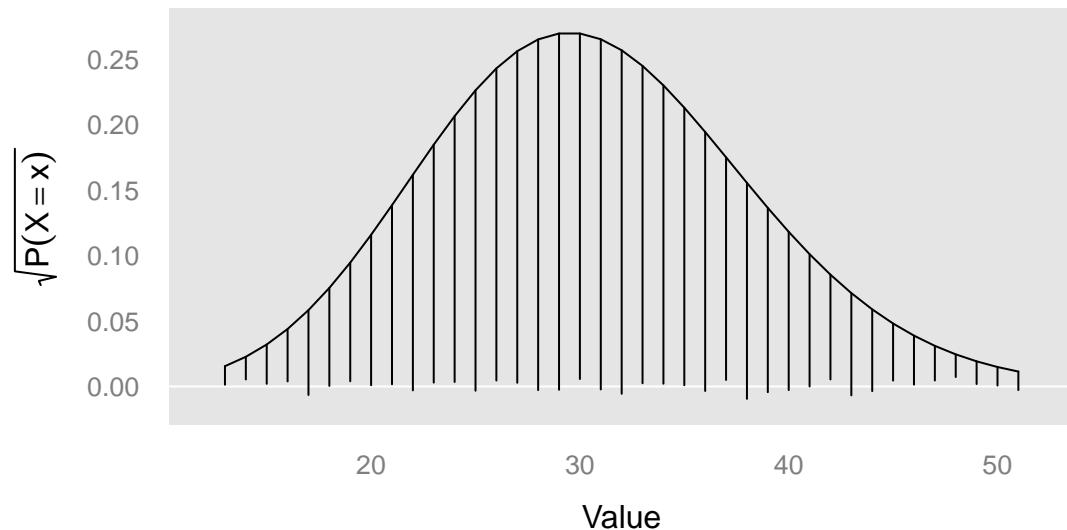


Figure 7.3: Histograms of random samples from the $\text{Normal}(0, 1)$ (top) and $\text{Exponential}(1)$ (bottom) distributions. Columns represent increasing sample sizes, and rows are independent repetitions of the experiment.

Samples from Poisson(30)



Samples from Binomial(60, 0.5)

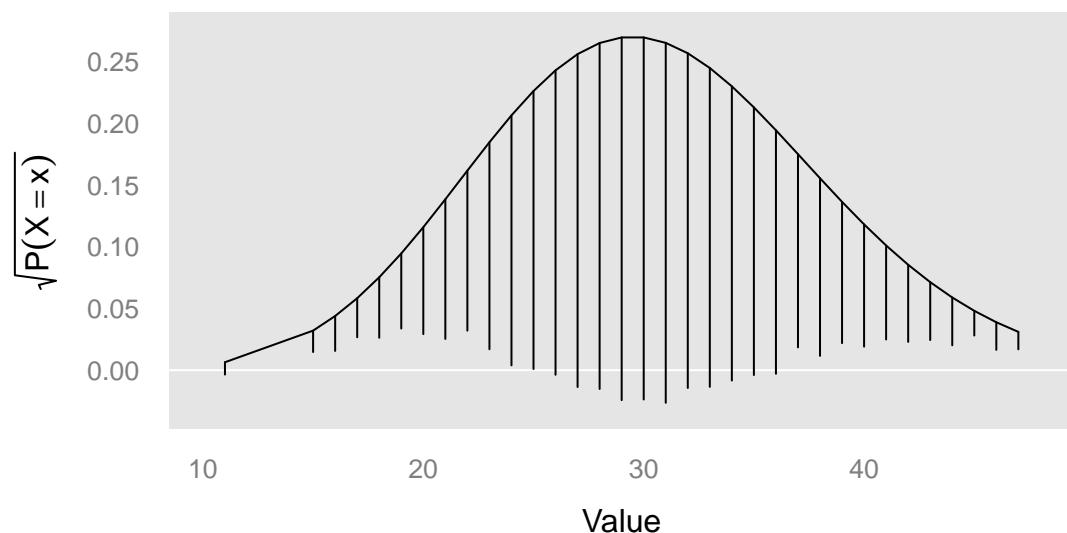


Figure 7.4: Hanging rootogram of 10000 random samples compared with the Poisson(30) distribution. In the top plot, the samples are also from Poisson(30), whereas in the bottom plot the samples are from the Binomial(100, 0.3) distribution, which has the same mean but different variance. Note the similarities with Figure 2.2

7.4.4 Q-Q Plots for Continuous Distributions

Just as histograms were binned versions of the empirical distribution plot, we can plot binned versions of hanging rootograms for data from a continuous distribution as well. It is more common however, to look at quantile-quantile plots (QQ plots), which do not bin the data, but instead plot what is essentially a transformation of the empirical CDF.

Recall that the ECDF of observations X_1, X_2, \dots, X_n is given by

$$\hat{F}_n(t) = P(Y \leq t) = \frac{\#\{X_i \leq t\}}{n}$$

The top plot in Figure 7.5 is a conventional ECDF plot of 200 observations simulated from a $\text{Normal}(1, 0.5^2)$ distribution. The bottom plot has the sorted data values on the y-axis and 200 equally spaced numbers from 0 to 1. A little thought tells us that this plot is essentially the same as the ECDF plot, with the x- and y-axes switched, and using points instead of lines. Naturally, we expect that for reasonably large sample sizes, the ECDF plot obtained from a random sample will be close to the true cumulative distribution function of the underlying distribution. If we know the shape of the distribution we expect the data to be from, we can compare it with the shape seen in the plot.

Although this is a fine idea in principle, it is difficult in practice to detect small differences between the observed shape and the theorized or expected shape. Here, we are helped again by the insight that the human eye finds it easier to detect deviations from a straight line than from curves. By keeping the sorted data values unchanged, but transforming the equally spaced probability values to the corresponding quantile values of the theorized distribution, we obtain a plot that we expect to be linear. Quantiles are defined as follows: For a given CDF F , the quantile corresponding to a probability value $p \in [0, 1]$ is a value x such that $F(x) = p$. Such an x may not exist for all p and F , and the definition of quantile needs to be modified to take this into account. However, for most standard continuous distributions used in Q-Q plots, the above definition is adequate. Such a plot with $\text{Normal}(0, 1)$ quantiles is shown in Figure 7.6.

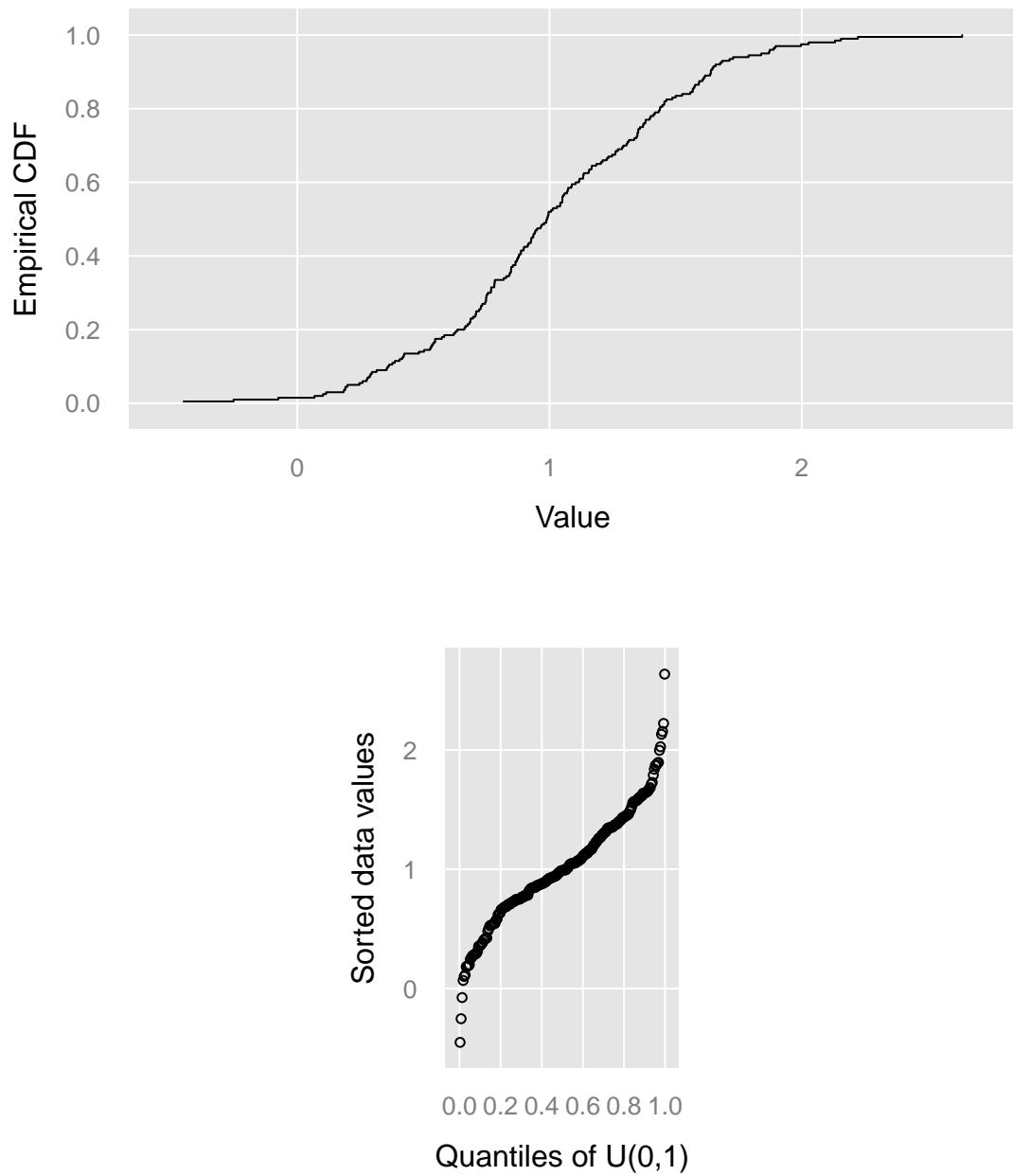


Figure 7.5: Conventional ECDF plot (top) and its “inverted” version (bottom), with x- and y-axes switched, and points instead of lines.

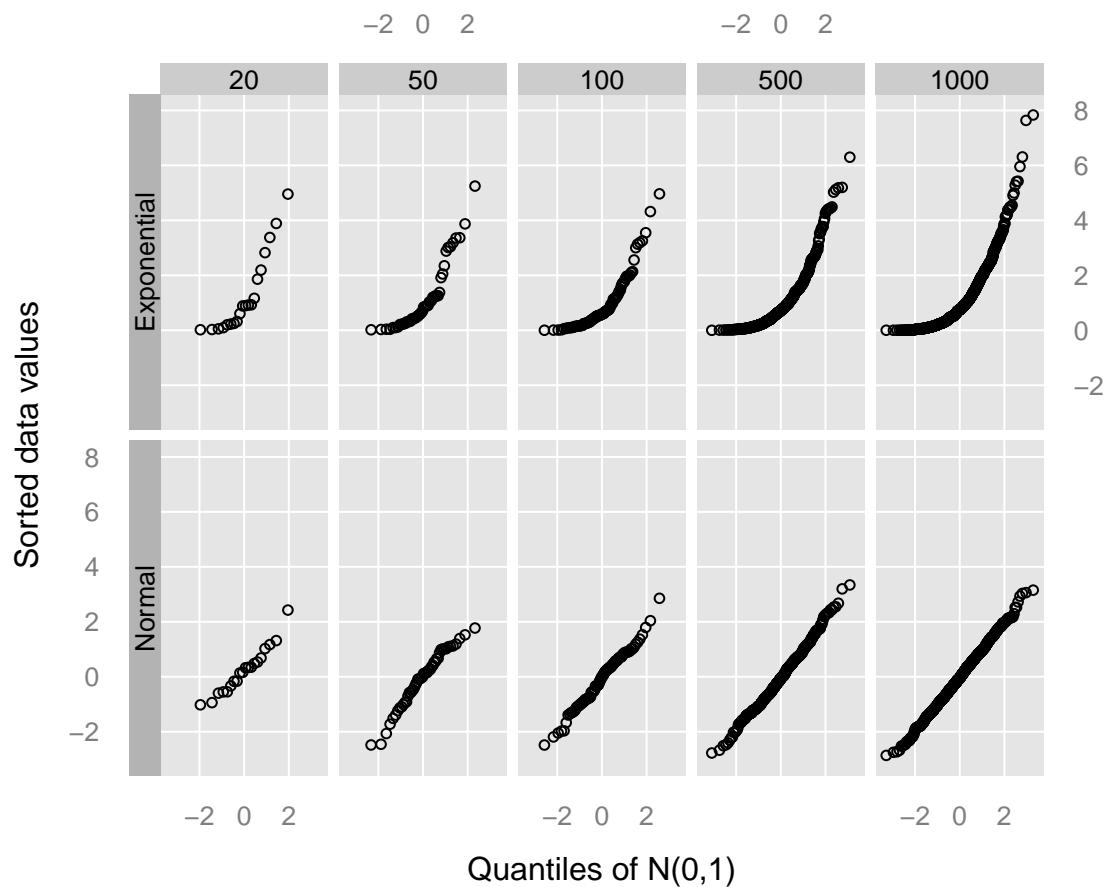


Figure 7.6: Normal Q-Q plots of data generated from Normal and Exponential distributions, with varying sample size. The Q-Q plots are more or less linear for Normal data, but exhibit curvature indicative of a relatively heavy right tail for exponential data. Not surprisingly, the difference becomes easier to see as the sample size increases.

8

SAMPLING DISTRIBUTIONS AND LIMIT THEOREMS

Let $n \geq 1$, X_1, X_2, \dots, X_n be an i.i.d. random sample from a population. Recall the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

and sample variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

We have seen in the previous chapter the significance of the above two sample statistics. In this chapter we shall discuss their distributional properties and limiting behaviour. In the next chapter we shall discuss how these results can be effectively used to verify specific hypotheses about the population. The corresponding field of study is called Hypothesis Testing or Test of Significance.

We will find the distribution of the sample mean and sample variance given the distribution of X_1 . One immediately observes that these are somewhat complicated functions of independent random variables. However in Section 3.3 and Section 5.5 we have seen examples of functions for which we were able to explicitly compute their distribution. To understand sampling statistics we must also understand the notion of joint distribution of more than two continuous random variables (See Section 3.3 for discrete random variables).

8.1 MULTI-DIMENSIONAL CONTINUOUS RANDOM VARIABLES

In Chapter 3, while discussing discrete random variables we had considered a finite collection of random variables (X_1, X_2, \dots, X_n) . In Definition 3.2.7, we had described how to define their joint distribution and we used this to understand the multinomial distribution in Example 3.2.12.

In the continuous setting as well there are many instances where it is relevant to study the joint distribution of a finite collection of random variables. Suppose X is a point chosen randomly in the unit sphere in the 3 dimensions. Then X has three coordinates and say $X = (X_1, X_2, X_3)$ where each X_i is a random variable in $(0, 1)$. Also they are dependent because we know that, $\sqrt{X_1^2 + X_2^2 + X_3^2} \leq 1$. It is useful and needed to understand their “joint distribution”. We have already seen the usefulness of sample mean and sample variance which are a function of X_1, X_2, \dots, X_n . To understand the distribution of sample mean and sample variance the joint distribution of X_1, X_2, \dots, X_n will be needed to be understood first. We define the joint distribution function first.

DEFINITION 8.1.1. Let $n \geq 1$ and X_1, X_2, \dots, X_n be random variables defined on the same probability space. The joint distribution function $F : \mathbb{R}^n \rightarrow [0, 1]$ is given by

$$F(x_1, x_2, \dots, x_n) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n),$$

for $x_1, x_2, \dots, x_n \in \mathbb{R}$.

As in single variable and two variable situations, the joint distribution function determines the entire joint distribution of X_1, X_2, \dots, X_n . More precisely, if all the random variables were discrete with $X_i : S \rightarrow T_i$ with T_i being countable subsets of $\subset \mathbb{R}$ for $1 \leq i \leq n$ from the joint distribution function one can determine

$$P(X_1 = t_1, X_2 = t_2, \dots, X_n = t_n),$$

for all $t_i \in T_i, 1 \leq i \leq n$. To understand the random variables in the continuous setting we need to set up some notation.

Let $n \geq 1$ and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a non-negative function, piecewise-continuous in each variable for which

$$\int_{\mathbb{R}^n} f(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n = 1.$$

For a Borel set $A \subset \mathbb{R}^n$ if

$$P(A) = \int_A f(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n.$$

Then one can show as in Theorem 5.1.5 that P is a probability on \mathbb{R}^n . f is called the density function for P . A sequence of random variables $(X_1, X_2, X_3, \dots, X_n)$ is said to have a joint density $f : \mathbb{R}^n \rightarrow \mathbb{R}$ if for every event $A \subset \mathbb{R}^n$

$$P((X_1, X_2, X_3, \dots, X_n) \in A) = \int_A f(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n.$$

In this setting, the joint distribution of (X_1, X_2, \dots, X_n) is determined by joint density f . Using multivariable calculus we can state and prove a similar type of result as Theorem 5.2.5 for random variables (X_1, X_2, \dots, X_n) that have a joint density. In particular, we can conclude that since the joint densities are assumed to be piecewise continuous, the corresponding distribution functions are piecewise differentiable. Further, the joint distribution of the continuous random variables (X_1, X_2, \dots, X_n) are completely determined by their joint distribution function F . That is, if we know $F(x_1, x_2, \dots, x_n)$ for all $x_1, x_2, \dots, x_n \in \mathbb{R}$ we could use multivariable calculus to differentiate F to find f . Then integrating this joint density over the event A we can calculate $P((X_1, X_2, \dots, X_n) \in A)$.

As in the $n = 2$ case one can recover the marginal density of each X_i for i between 1 and n by integrating over the other indices. So, the marginal density of X_i at a is given by

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

Further for $n \geq 3$, we can deduce the joint density for any sub-collection $m \leq n$ random variables by integrating over the other variables. For instance, if we were interested in the joint density of (X_1, X_3, X_7) we would obtain

$$f_{X_1, X_3, X_7}(a_1, a_3, a_7) = \int_{\mathbb{R}^{n-3}} f(a_1, x_2, a_3, x_4, x_5, x_6, a_7, x_8 \dots, x_n) dx_2 dx_4 dx_5 dx_6 \dots dx_n.$$

Suppose X_1, X_2, \dots, X_n are random variables defined on a single sample space S with joint density $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function of n variables for which $g(X_1, X_2, \dots, X_n)$ is defined in the range of the X_j variables. Let B be an event in the range of g . Then, following the proof of Theorem 3.3.5, we can show that

$$P(g(X_1, X_2, \dots, X_n) \in B) = P((X_1, X_2, \dots, X_n) \in g^{-1}(B)).$$

The above provides an abstract method of finding the distribution of the random variable $Y = g(X_1, X_2, \dots, X_n)$ but it might be difficult to calculate it explicitly. For $n = 1$, we discussed this question in detail in Section 5.3, for $n = 2$ we did explore how to find the distributions of sums and ratios of independent random variables (see Section 5.5). In a few cases by induction n , this method could be extended but in general it is not possible. In Appendix B, Section B.2 we discuss the Jacobian method of finding the joint density of the transformed random variable.

The notion of independence also extends to the multi-dimensional continuous random variable as in discrete setting. As discussed in Definition 3.2.3, a finite collection of continuous random variables X_1, X_2, \dots, X_n is mutually independent if the sets $(X_j \in A_j)$ are mutually independent for all events A_j in the ranges of the corresponding X_j . As proved for the $n = 2$ case in Theorem 5.4.7, we can similarly

deduce that if $(X_1, X_2, X_3, \dots, X_n)$ are mutually independent continuous random variables with marginal densities f_{X_i} then their joint density is given by

$$f(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i), \quad (8.1.1)$$

for $x_i \in \mathbb{R}$ and $1 \leq i \leq n$. Further for any finite sub-collection $\{X_{i_1}, X_{i_2}, \dots, X_{i_m}\}$ of the above independent random variables, their joint density is given by

$$f(a_1, a_2, \dots, a_m) = \prod_{i=j}^m f_{X_{i_j}}(a_j). \quad (8.1.2)$$

We conclude this section with a result that we will repeatedly use.

THEOREM 8.1.2. *For each $j \in \{1, 2, \dots, n\}$ define a positive integer m_j and suppose $X_{i,j}$ is an array of mutually independent continuous random variables for $j \in \{1, 2, \dots, n\}$ and $i \in \{1, 2, \dots, m_j\}$. Let $g_j(\cdot)$ be functions such that the quantity*

$$Y_j = g_j(X_{1,j}, X_{2,j}, \dots, X_{m_j,j})$$

is defined for the outputs of the $X_{i,j}$ variables. Then the resulting variables Y_1, Y_2, \dots, Y_n are mutually independent.

Proof- Follows by the same proof presented in Theorem 3.3.6. ■

8.1.1 Order Statistics and their Distributions

Let $n \geq 1$ and let X_1, X_2, \dots, X_n be a i.i.d random sample from a population. Let F be the common distribution function. Let the X 's be arranged in increasing order of magnitude denoted by

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}.$$

These ordered values are called the order statistics of the sample X_1, X_2, \dots, X_n . For, $1 \leq r \leq n$, $X_{(r)}$ is called the r -th order statistic. One can computer $F_{(r)}$, the distribution function of $X_{(r)}$, for $1 \leq r \leq n$ in terms of n and F . We have,

$$\begin{aligned} F_{(1)}(x) &= P(X_{(1)} \leq x) = 1 - P(X_{(1)} > x) = 1 - P(\cap_{i=1}^n (X_i > x)) \\ &= 1 - \prod_{i=1}^n P(X_i > x) = 1 - \prod_{i=1}^n (1 - P(X_i \leq x)) \\ &= 1 - (1 - F(x))^n, \end{aligned}$$

$$F_{(n)}(x) = P(X_{(n)} \leq x) = P(\cap_{i=1}^n (X_i \leq x)) = \prod_{i=1}^n P(X_i \leq x) = (F(x))^n,$$

and for $1 < r < n$,

$$\begin{aligned} F_{(r)}(x) &= P(X_{(r)} \leq x) = P(\text{at least } r \text{ elements from the sample are } \leq x) \\ &= \sum_{j=r}^n P(\text{exactly } j \text{ elements from the sample are } \leq x) \\ &= \sum_{j=r}^n \binom{n}{j} P(\text{chosen } j \text{ elements from the sample are } \leq x) \times \\ &\quad \times P((n-j) \text{ elements not chosen from the sample are } > x) \\ &= \sum_{j=r}^n \binom{n}{j} F(x)^j (1 - F(x))^{n-j} \end{aligned}$$

If the distribution function F had a probability density function f then each $X_{(r)}$ has a probability density function $f_{(r)}$. This can be obtained by differentiating $F_{(r)}$ and is given by the below expression.

$$f_{(r)}(x) = \begin{cases} n(1 - F(x))^{n-1} f(x) & r = 1 \\ nf(x)(F(x))^{n-1} & r = n \\ \frac{n!}{(r-1)!(n-r)!} f(x)(F(x))^{r-1}(1 - F(x))^{n-r} & 1 < r < n \end{cases} \quad (8.1.3)$$

EXAMPLE 8.1.3. Let $n \geq 1$ and let X_1, X_2, \dots, X_n be a i.i.d random sample from a population whose common distribution F is an Exponential (λ) random variable. Then we know that

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

Therefore using (8.1.3) and substituting for F as above we have that the densities of the order statistics are given by

$$f_{(r)}(x) = \begin{cases} n(e^{-\lambda x})^{n-1} \lambda e^{-\lambda x} & r = 1 \\ n\lambda e^{-\lambda x}(1 - e^{-\lambda x})^{n-1} & r = n \\ \lambda e^{-\lambda x} \frac{n!}{(r-1)!(n-r)!} (1 - e^{-\lambda x})^{r-1} (e^{-\lambda x})^{n-r} & 1 < r < n, \end{cases}$$

for $x > 0$. Simplifying the algebra we obtain,

$$f_{(r)}(x) = \begin{cases} n\lambda e^{-n\lambda x} & r = 1 \\ n\lambda e^{-\lambda x}(1 - e^{-\lambda x})^{n-1} & r = n \\ \frac{\lambda n!}{(r-1)!(n-r)!} (1 - e^{-\lambda x})^{r-1} (e^{-\lambda x})^{n-r+1} & 1 < r < n, \end{cases}$$

for $x > 0$. We note from the above that $X_{(1)}$, i.e minimum of exponentials, is Exponential ($n\lambda$) random variable. However the other order statistics are not exponentially distributed. ■

In many applications one is interested in the range of values a random variable X assumes. A method to understand this is to sample X_1, X_2, \dots, X_n i.i.d. X and examine $R = X_{(n)} - X_{(1)}$. Suppose X has a probability density function $f : \mathbb{R} \rightarrow \mathbb{R}$ and distribution function $F : \mathbb{R} \rightarrow [0, 1]$. As before we can calculate the joint density of $X_{(1)}, X_{(n)}$ by first computing the joint distribution function. This is done by using the i.i.d. nature of the sample and the definition of the order statistics.

$$\begin{aligned} P(X_{(1)} \leq x, X_{(n)} \leq y) &= P(X_{(n)} \leq y) - P(x < X_{(1)}, X_{(n)} \leq y) \\ &= P(\cap_{i=1}^n \{X_i \leq y\}) - P(\cap_{i=1}^n \{x < X_i \leq y\}) \\ &= [P(X \leq y)]^n - [P(x < X \leq y)]^n \\ &= \begin{cases} [F(x)]^n - [F(y) - F(x)]^n & x < y \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

From the above, differentiating partially in x and y we see that the joint density of $(X_{(1)}, X_{(n)})$ is given by

$$f_{X_{(1)}, X_{(n)}}(x, y) = \begin{cases} n(f(x) - f(y))[F(y) - F(x)]^{n-1} & x < y \\ 0 & \text{otherwise.} \end{cases} \quad (8.1.4)$$

To calculate the distribution of R , we compute its distribution function. For $r \leq 0$, $P(R \leq r) = 0$ and for $r > 0$, using the above joint density of $(X_{(1)}, X_{(n)})$ we have

$$\begin{aligned} P(R \leq r) &= P(X_{(n)} \leq X(1) + r) \\ &= \int_{-\infty}^{\infty} \left[\int_0^r f_{X_{(1)}, X_{(n)}}(x, z+x) dz \right] dx \\ &= \int_0^r \left[\int_{-\infty}^{\infty} f_{X_{(1)}, X_{(n)}}(x, z+x) dz \right] dx, \end{aligned}$$

where we have done a change of variable $y = z + x$ in the second last line and a change in the order of integration in the last line. Differentiating the above we conclude that R has a joint density given by

$$f_R(r) = \begin{cases} \int_{-\infty}^{\infty} f_{X_{(1)}, X_{(n)}}(x, r+x) dx & \text{if } r > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (8.1.5)$$

EXAMPLE 8.1.4. Let X_1, X_2, \dots, X_n be i.i.d Uniform $(0, 1)$. The probability density function and distribution function of a uniform $(0, 1)$ random variable are given by

$$f(x) = \begin{cases} 1 & \text{if } x \in (0, 1) \\ 0 & \text{otherwise.} \end{cases} \quad \text{and} \quad F(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } 0 < x < 1 \\ 1 & \text{if } x > 1. \end{cases}$$

Using (8.1.3), we have the probability density function of

$$\begin{aligned} X_{(1)} \text{ is given by} \quad f_{X_{(1)}}(x) &= \begin{cases} n(1-x)^{n-1} & \text{if } x \in (0, 1) \\ 0 & \text{otherwise,} \end{cases} \\ X_{(n)} \text{ is given by} \quad f_{X_{(n)}}(x) &= \begin{cases} nx^{n-1} & \text{if } x \in (0, 1) \\ 0 & \text{otherwise.} \end{cases} \\ X_{(r)} \text{ is given by} \quad f_{X_{(r)}}(x) &= \begin{cases} \frac{n!}{(r-1)!(n-r)!} x^{r-1} (1-x)^{n-r} & \text{if } x \in (0, 1) \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

for $1 < r < n$.

Using (8.1.4), we have the joint density of

$$(X_{(1)}, X_{(n)}) \text{ is given by} \quad f_{X_{(1)}, X_{(n)}}(x, y) = \begin{cases} n(n-1)(y-x)^{n-1} & \text{if } 0 \leq x \leq y \leq 1 \\ 0 & \text{otherwise,} \end{cases}$$

Using (8.1.5), we have the probability density function of the range

$$\begin{aligned} R = X_{(n)} - X_{(1)} \text{ is given by } f_R(r) &= \begin{cases} \int_0^{1-r} n(n-1)(x+r-x)^{n-1} dx & \text{if } 0 < r < 1 \\ 0 & \text{otherwise,} \end{cases} \\ &= \begin{cases} n(n-1)r^{n-1}(1-r) & \text{if } 0 < r < 1 \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

We see that $X_{(r)} \sim \text{Beta}(r, n-r+1)$ for $1 \leq r \leq n$ and the range $R \sim \text{Beta}(n, 2)$



In general we can also understand the joint-distribution of the order statistics. Suppose we have an i.i.d sample X_1, X_2, \dots, X_n having distribution X . If X has a probability density function $f : \mathbb{R} \rightarrow \mathbb{R}$ then one can show that the order statistic $(X_{(1)}, X_{(2)}, \dots, X_{(n)})$ has a joint density $h : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$h(u_1, u_2, \dots, u_n) = \begin{cases} n!f(u_1)f(u_2)\dots f(u_n) & u_1 < u_2 < \dots < u_n, \\ 0 & \text{otherwise.} \end{cases}$$

The above fact intuitively is clear. Any ordering $u_1 < u_2 < \dots < u_n$ "has a probability" $f(u_1)f(u_2)\dots f(u_n)$. Each of the X_i can assume any of the u_k 's. The total number of possible orderings is $n!$. A formal proof involves using the Jacobian method and will be discussed in Appendix B.

8.1.2 χ^2 , F and t

χ^2 , F and t distributions arise naturally when considering functions of i.i.d. normal random variables $(X_1, X_2, X_3, \dots, X_n)$ for $n \geq 1$. They also are useful in Hypothesis testing as well. We discuss these via three examples.

EXAMPLE 8.1.5. (Chi-Square) Let $n \geq 1$ and $(X_1, X_2, X_3, \dots, X_n)$ be a collection of independent Normal random variables with mean 0 and variance 1. Then the joint density is given by

$$f(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i) = \frac{1}{\sqrt[4]{2\pi}} e^{-\sum_{i=1}^n \frac{x_i^2}{2}},$$

for $x_i \in \mathbb{R}$ and $1 \leq i \leq n$.

Let $Z = \sum_{i=1}^n X_i^2$. We shall find the distribution of Z in two steps. First, clearly the range of X_1^2 is non-negative. The distribution function for X_1^2 at $z \geq 0$, is given by

$$\begin{aligned} F_1(z) &= P(X_1^2 \leq z) \\ &= P(X_1 \leq \sqrt{z}) \\ &= \int_0^{\sqrt{z}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \\ &= \int_0^z \frac{1}{2\sqrt{2\pi}} e^{-\frac{u}{2}} u^{-\frac{1}{2}} du \end{aligned}$$

Comparing it with the Gamma (α, λ) random variable defined in Definition 5.5.5 and using Exercise 5.5.10, we see that X_1^2 is distributed as a Gamma $(\frac{1}{2}, \frac{1}{2})$ random variable. Using the calculation done in Example 5.5.6 for $n = 2$ and by induction we have that $Z = \sum_{i=1}^n X_i^2$ will be Gamma $(\frac{n}{2}, \frac{1}{2})$. This distribution is referred to as Chi-Square with n - degrees of freedom. We define it precisely next. ■

DEFINITION 8.1.6. (Chi-Square with n degrees of freedom) A random variable X whose distribution is Gamma $(\frac{n}{2}, \frac{1}{2})$ is said to have Chi-square distribution with n -degrees of freedom (i.e number of parameters). Gamma $(\frac{n}{2}, \frac{1}{2})$ is denoted by χ_n^2 and as discussed earlier it has density given by

$$\begin{aligned} f(x) &= \frac{2^{-\frac{n}{2}}}{\Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} \\ &= \begin{cases} \frac{2^{-\frac{n}{2}}}{(\frac{n}{2}-1)!} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} & \text{when } n \text{ is even.} \\ \frac{2^{n-\frac{n}{2}-1} (\frac{n-1}{2})!}{(n-1)! \sqrt{\pi}} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} & \text{when } n \text{ is odd.} \end{cases} \end{aligned}$$

when $x > 0$.

We shall show in the next subsection that sample variance from a Normal population is a Chi-square random variable. In the next chapter we shall construct a test to make inferences about the variances of the two population. In that context we shall compare sample variances and this is where the F distribution arises naturally.

EXAMPLE 8.1.7. (F-distribution) Suppose X_1, X_2, \dots, X_{n_1} be an i.i.d. random sample from a Normal mean 0 and variance σ_1^2 population and Y_1, Y_2, \dots, Y_{n_2} be an i.i.d. random sample from a Normal mean 0 and variance σ_2^2 population. We have already seen in Example 8.1.5 that $U = \sum_{i=1}^{n_1} \left(\frac{X_i}{\sigma_1}\right)^2$ is a $\chi_{n_1}^2$ random variable and $V = \sum_{i=1}^{n_2} \left(\frac{Y_i}{\sigma_2}\right)^2$ is a $\chi_{n_2}^2$ random variable. Further U and V are independent. Let $Z = \frac{U}{n_1} / \frac{V}{n_2}$. Let $Y = \frac{n_1}{n_2} Z = \frac{U}{V}$. As done in Example 5.5.10 the density of Y for $y > 0$ is given by

$$f_Y(y) = \frac{y^{\frac{n_1}{2}-1}}{(1+y)^{\frac{n_1+n_2}{2}}} \frac{\Gamma(\frac{n_1+n_2}{2})}{\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})}$$

Therefore, for $z > 0$

$$\begin{aligned} F_Z(z) &= P(Z \leq z) \\ &= P(Y \leq \frac{n_2}{n_1} z) \\ &= \int_{-\infty}^{\frac{n_2}{n_1} z} \frac{y^{\frac{n_1}{2}-1}}{(1+y)^{\frac{n_1+n_2}{2}}} \frac{\Gamma(\frac{n_1+n_2}{2})}{\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} dy \\ &\quad \text{making a } u\text{-substitution with } \frac{n_1}{n_2} y = u \\ &= \int_{-\infty}^z \left(\frac{n_2}{n_1}\right)^{\frac{n_1}{2}} \frac{u^{\frac{n_1}{2}-1}}{(1+\frac{n_1}{n_2}u)^{\frac{n_1+n_2}{2}}} \frac{\Gamma(\frac{n_1+n_2}{2})}{\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})} du \end{aligned}$$

Therefore the density of Z , for $z > 0$ is given by

$$f(z) = \left(\frac{n_2}{n_1}\right)^{\frac{n_1}{2}} \frac{z^{\frac{n_1}{2}-1}}{(1+\frac{n_1}{n_2}z)^{\frac{n_1+n_2}{2}}} \frac{\Gamma(\frac{n_1+n_2}{2})}{\Gamma(\frac{n_1}{2})\Gamma(\frac{n_2}{2})}.$$

Z is said to have $F(n_1, n_2)$ distribution. Z is close to a widely used distribution in statistics called F -distribution. ■

The distribution of the ratio of sample mean and sample variance plays an important role in Hypothesis testing. This forms the motivation for the next example where the t distribution arises naturally.

EXAMPLE 8.1.8. (t-distribution) Let X_1 be a Normal random variable with mean 0 and variance 1. Let X_2 be an independent χ_n^2 random variable. Let

$$Z = \frac{X_1}{\sqrt{\frac{X_2}{n}}}.$$

We wish to find the density of Z . Observe that $U = Z^2$ is given by $\frac{X_1^2}{\frac{X_2}{n}}$. Now, X_1^2 has χ_1^2 distribution (See Example 8.1.5). So applying Example 8.1.7 with $n_1 = 1$ and $n_2 = n$, we find that U has $F(1, n)$ distribution. The density of U is given by

$$\begin{aligned} f_U(u) &= \left(\frac{1}{n}\right)^{\frac{1}{2}} \frac{u^{\frac{1}{2}-1}}{(1+\frac{1}{n}u)^{\frac{n+1}{2}}} \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{n}{2})} \\ &= \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \frac{u^{-\frac{1}{2}}}{(1+\frac{u}{n})^{\frac{n+1}{2}}}. \end{aligned}$$

Since X_1 is a symmetric random variable and $\sqrt{\frac{X_2}{n}}$ is positive valued we conclude that Z is a symmetric random variable (Exercise 8.1.10). So, for $u > 0$

$$\begin{aligned} P(U \leq u) &= P(Z^2 \leq u) \\ &= P(-\sqrt{u} \leq Z \leq \sqrt{u}) \\ &= P(Z \leq \sqrt{u}) - P(Z \leq -\sqrt{u}) \\ &= P(Z \leq \sqrt{u}) - P(Z \geq \sqrt{u}) \\ &= 2P(Z \leq \sqrt{u}) - 1 \end{aligned}$$

Therefore if $f_Z(\cdot)$ is the density of Z then

$$f_U(u) = \frac{1}{\sqrt{u}}(f_Z(\sqrt{u})).$$

Hence for any $z \in \mathbb{R}$ the density of Z is given by

$$\begin{aligned} f_Z(z) &= |z| f_U(z^2) \\ &= |z| \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \frac{z^{2-\frac{1}{2}}}{(1+\frac{u}{n})^{\frac{n+1}{2}}} \\ &= \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \left(1 + \frac{z^2}{n}\right)^{-\frac{n+1}{2}} \end{aligned}$$

Z is said to have t -distribution with n -degrees of freedom. We will denote this by the notation $Z \sim t_n$. ■

8.1.3 Distribution of Sampling Statistics from a Normal population

Let $n \geq 1$, X_1, X_2, \dots, X_n be an i.i.d. random sample from a population having mean μ and variance σ^2 . Consider the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

and sample variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

We have already seen in Theorem 7.2.2 that $E[\bar{X}] = \mu$ and in Theorem 7.2.4 that $E[S^2] = \sigma^2$. It turns out that it is not easy to understand the precise distribution of \bar{X} or S^2 in general. However, this can be done when the population is normally distributed. The main result of this section is the following.

THEOREM 8.1.9. *Let $n \geq 1$, X_1, X_2, \dots, X_n , be an i.i.d random sample with distribution $X \sim \text{Normal}(\mu, \sigma^2)$. Let \bar{X} and S^2 be as above. Then,*

- (a) \bar{X} is a Normal random variable with mean μ and variance $\frac{\sigma^2}{n}$.
- (b) $\frac{(n-1)}{\sigma^2} S^2$ has the χ_{n-1}^2 distribution.
- (c) \bar{X} and S^2 are independent.

Proof - (a) follows from Theorem 6.3.13. The typical proof for (b) and (c) is via Helmert's transformation (see [Rao73]) and requires some knowledge of Linear Algebra. Here we will follow Kruskal's proof as illustrated in [Stig84]. The proof is by the method of induction. For implementing the inductive step on the sample size n , we shall replace \bar{X} and S^2 with \bar{X}_n and S_n^2 for the rest of the proof.

Step 1: (Proof for $n = 2$) Here

$$\bar{X}_2 = \frac{X_1 + X_2}{2} \quad \text{and} \quad S_2^2 = \left(X_1 - \frac{X_1 + X_2}{2} \right)^2 + \left(X_2 - \frac{X_1 + X_2}{2} \right)^2 = \frac{(X_1 - X_2)^2}{2}. \quad (8.1.6)$$

(a) Follows from Theorem 6.3.13.

(b) As X_1 and X_2 are independent Normal random variables with mean μ and variance σ^2 , by Theorem 6.3.13, $\frac{(X_1 - X_2)}{\sigma\sqrt{2}}$ is a Normal random variable with mean 0 and variance 1. Using Example 8.1.5, we know that $\frac{S_2^2}{\sigma^2}$ has χ_1^2 distribution and this proves (b).

(c) From (8.1.6), \bar{X}_2 is a function of $X_1 + X_2$ and S_2^2 is a function of $X_1 - X_2$. Theorem 8.1.2 will imply that \bar{X}_2 and S_2^2 are independent if we show $X_1 + X_2$ and $X_1 - X_2$ are independent. Let $\alpha, \beta \in \mathbb{R}$. Then using Theorem 6.3.13 again we have that $\alpha(X_1 + X_2) + \beta(X_1 - X_2) = (\alpha + \beta)X_1 + (\alpha - \beta)X_2$ is normally distributed. As this is true for any $\alpha, \beta \in \mathbb{R}$, by Definition 6.4.1 $(X_1 + X_2, X_1 - X_2)$ is a bivariate normal random variable. Using Theorem 6.2.2 (f) and (g), along with the fact that X_1 and X_2 are independent Normal random variables with mean μ and variance σ^2 , we have

$$\text{Cov}[X_1 + X_2, X_1 - X_2] = \text{Var}[X_1] + \text{Cov}[X_2, X_1] - \text{Cov}[X_1, X_2] - \text{Var}[X_2] = 0.$$

Theorem 6.4.3 then implies that $X_1 + X_2$ and $X_1 - X_2$ are independent.

Step 2: (inductive hypothesis) Let us inductively assume that (a), (b), and (c) are true when $n = k$ for some $k \in \mathbb{N}$.

Step 3: (Proof for $n = k + 1$) We shall rewrite \bar{X}_{k+1} and S_{k+1}^2 using some elementary algebra.

$$\bar{X}_k - \bar{X}_{k+1} = \bar{X}_k - \frac{1}{k+1} \sum_{i=1}^{k+1} X_i = \left(1 - \frac{k}{k+1}\right) \bar{X}_k - \frac{1}{k+1} X_{k+1} = \frac{1}{k+1} (\bar{X}_k - X_{k+1}). \quad (8.1.7)$$

Adding and subtracting \bar{X}_k inside the summand of S_{k+1}^2 , we have

$$\begin{aligned} S_{k+1}^2 &= \frac{1}{k} \sum_{i=1}^{k+1} (X_i - \bar{X}_{k+1})^2 = \frac{1}{k} \sum_{i=1}^{k+1} (X_i - \bar{X}_k + \bar{X}_k - \bar{X}_{k+1})^2 \\ &= \frac{1}{k} \sum_{i=1}^{k+1} (X_i - \bar{X}_k)^2 + 2(X_i - \bar{X}_k)(\bar{X}_k - \bar{X}_{k+1}) + (\bar{X}_k - \bar{X}_{k+1})^2 \\ &= \frac{k-1}{k} S_k^2 + \frac{1}{k} (X_{k+1} - \bar{X}_k)^2 + \frac{1}{k} (2(X_{k+1} - \bar{X}_k)(\bar{X}_k - \bar{X}_{k+1}) + (k+1)(\bar{X}_k - \bar{X}_{k+1})^2) \\ &= \frac{k-1}{k} S_k^2 + \frac{1}{k} (X_{k+1} - \bar{X}_k)^2 - \frac{1}{k} \left(2(X_{k+1} - \bar{X}_k) \frac{(X_{k+1} - \bar{X}_k)}{k+1} + \frac{(X_{k+1} - \bar{X}_k)^2}{k+1} \right) \\ &= \frac{k-1}{k} S_k^2 + \frac{1}{k+1} (X_{k+1} - \bar{X}_k)^2, \end{aligned}$$

where we have used (8.1.7) in the second last inequality. Dividing across by σ^2 and multiplying by k we have

$$\frac{k}{\sigma^2} S_{k+1}^2 = \frac{k-1}{\sigma^2} S_k^2 + \frac{k}{\sigma^2(k+1)} (X_{k+1} - \bar{X}_k)^2. \quad (8.1.8)$$

(a) Follows from Theorem 6.3.13.

(b) To prove (b), it is enough to show that:

$\left(\sqrt{\frac{k}{(k+1)\sigma^2}} \right) (X_{k+1} - \bar{X}_k)$ is a standard normal random variable and is independent of $\frac{(k-1)}{\sigma^2} S_k^2$.

The reason being: $\frac{k}{\sigma^2(k+1)} (X_{k+1} - \bar{X}_k)^2$ then has χ_1^2 distribution by Example 8.1.5 and is independent of $\frac{(k-1)}{\sigma^2} S_k^2$ by Theorem 8.1.2; by the induction hypothesis $\frac{(k-1)}{\sigma^2} S_k^2$ has the χ_{k-1}^2 distribution; and finally using (8.1.8) along with Example 5.5.6, will imply that $\frac{k}{\sigma^2} S_{k+1}^2$ has χ_k^2 distribution.

As

$$\left(\sqrt{\frac{k}{(k+1)\sigma^2}} \right) (X_{k+1} - \bar{X}_k) = \left(\sqrt{\frac{(k+1)\sigma^2}{k}} \right) X_{k+1} - \sum_{i=1}^k \frac{1}{k} \left(\sqrt{\frac{k}{(k+1)\sigma^2}} \right) X_i$$

It is routine calculation using Theorem 6.3.13 to see that is a standard normal random variable.

By induction hypothesis \bar{X}_k and $\frac{k-1}{\sigma^2} S_k^2$ are independent. Since X_1, \dots, X_k, X_{k+1} are mutually independent, Theorem 8.1.2 implies that X_{k+1} is independent of \bar{X}_k and $\frac{k-1}{\sigma^2} S_k^2$. Therefore,

$$\bar{X}_k, \quad \frac{k-1}{\sigma^2} S_k^2, \quad X_{k+1} \quad \text{are mutually independent random variables.} \quad (8.1.9)$$

Consequently, another application of Theorem 8.1.2 will then imply that $\frac{k}{\sigma^2(k+1)}(X_{k+1} - \bar{X}_k)^2$ and $\frac{(k-1)}{\sigma^2} S_k^2$ are independent random variables.

(c) To prove (c), it is enough to show that \bar{X}_{k+1} and $X_{k+1} - \bar{X}_k$ are independent. The reason is the following:

- (i) Theorem 8.1.2 then implies \bar{X}_{k+1} is independent of $\frac{k}{\sigma^2(k+1)}(X_{k+1} - \bar{X}_k)^2$;
- (ii) \bar{X}_{k+1} is a function of X_{k+1} and \bar{X}_k . So (8.1.9) and Theorem 8.1.2 will then imply \bar{X}_{k+1} is independent of $\frac{(k-1)}{\sigma^2} S_k^2$ and also $\frac{k}{\sigma^2(k+1)}(X_{k+1} - \bar{X}_k)^2$ is independent of $\frac{(k-1)}{\sigma^2} S_k^2$;
- (iii) Using (i) and (ii) we can conclude that \bar{X}_{k+1} , $\frac{(k-1)}{\sigma^2} S_k^2$, and $\frac{k}{\sigma^2(k+1)}(X_{k+1} - \bar{X}_k)^2$ are mutually independent; and
- (iv) finally S_{k+1}^2 is a function $\frac{(k-1)}{\sigma^2} S_k^2$, and $\frac{k}{\sigma^2(k+1)}(X_{k+1} - \bar{X}_k)^2$ by (8.1.8). Then (iii) and Theorem 8.1.2 will imply that S_{k+1}^2 and \bar{X}_{k+1} are independent.

Let $\alpha, \beta \in \mathbb{R}$. We have

$$\alpha(\bar{X}_{k+1}) + \beta(X_{k+1} - \bar{X}_k) = \sum_{i=1}^k \left(\frac{\alpha}{k+1} - \frac{\beta}{k} \right) X_i + \left(\frac{\alpha}{k+1} - \beta \right) X_{k+1}.$$

Theorem 6.3.13 will imply that $\alpha(\bar{X}_{k+1}) + \beta(X_{k+1} - \bar{X}_k)$ is a normally distributed random variable for any $\alpha, \beta \in \mathbb{R}$. So by Definition 6.4.1 $(\bar{X}_{k+1}, X_{k+1} - \bar{X}_k)$ is a bivariate normal random variable. Further, from Theorem 6.2.2 (f) and (g), we have

$$\begin{aligned} Cov[\bar{X}_{k+1}, X_{k+1} - \bar{X}_k] &= Cov[\frac{k\bar{X}_k + X_{k+1}}{k+1}, X_{k+1} - \bar{X}_k] \\ &= \frac{1}{k+1} Var[X_{k+1}] - Cov[\bar{X}_k, X_{k+1}] - \frac{k}{k+1} Var[\bar{X}_k] \\ &= \frac{1}{k+1} \sigma^2 + 0 - \frac{k}{k+1} \frac{\sigma^2}{k} = 0, \end{aligned}$$

where we have used (8.1.9) in the last line. From Theorem 6.4.3 we conclude that $\bar{X}_{k+1}, X_{k+1} - \bar{X}_k$ are independent. ■

The following Corollary will be used in Chapter 9

COROLLARY 8.1.10. Let $n \geq 1$, X_1, X_2, \dots, X_n , be an i.i.d random sample with distribution $X \sim Normal(\mu, \sigma^2)$. Let \bar{X} and S^2 be as above. Then

$$\frac{\sqrt{n}(\bar{X} - \mu)}{S}$$

has the t_{n-1} distribution.

Proof - From Theorem 8.1.9 it is clear that

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

is a Normal random variable with mean 0 and variance 1, and

$$\frac{(n-1)}{\sigma^2} S^2$$

is a χ_{n-1}^2 random variable. Note

$$\frac{\sqrt{n}(\bar{X} - \mu)}{S} = \frac{\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{1}{n-1} \frac{(n-1)S^2}{\sigma^2}}}$$

So by Example 8.1.8 we have the result. ■

EXERCISES

Ex. 8.1.1. Verify that each of $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ are density functions on \mathbb{R}^3 .

- (a) $f(x_1, x_2, x_3) = \begin{cases} \frac{2}{3}(x_1 + x_2 + x_3) & \text{if } 0 < x_i < 1, i = 1, 2, 3. \\ 0 & \text{otherwise} \end{cases}$
- (b) $f(x_1, x_2, x_3) = \begin{cases} \frac{1}{8}(x_1^2 + x_2^2 + x_3^2) & \text{if } 0 < x_i < 2, i = 1, 2, 3. \\ 0 & \text{otherwise} \end{cases}$
- (c) $f(x_1, x_2, x_3) = \begin{cases} \frac{2}{81}x_1x_2x_3 & \text{if } 0 < x_i < 3, i = 1, 2, 3. \\ 0 & \text{otherwise} \end{cases}$
- (d) $f(x_1, x_2, x_3) = \begin{cases} \frac{3}{4}(x_1x_2 + x_1x_3 + x_2x_3) & \text{if } 0 < x_i < 1, i = 1, 2, 3. \\ 0 & \text{otherwise} \end{cases}$

Ex. 8.1.2. Suppose (X_1, X_2, X_3) have a joint density $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ given by

$$f(x_1, x_2, x_3) = \begin{cases} \frac{4}{3}(x_1^3 + x_2^3 + x_3^3) & \text{if } 0 < x_i < 1, i = 1, 2, 3. \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find $P(X_1 < \frac{1}{2}, X_3 > \frac{1}{2})$.
- (b) Find the joint density of $(X_1, X_2), (X_1, X_3), (X_2, X_3)$.
- (c) Find the marginal densities of X_1, X_2 , and X_3 .

Ex. 8.1.3. Let D be a set in \mathbb{R}^3 with a well defined volume. (X_1, X_2, X_3) are said be uniform on a set D if they have a joint density given by

$$f(x_1, x_2, x_3) = \begin{cases} \frac{1}{\text{Volume}(D)} & \text{if } x \in D \\ 0 & \text{otherwise.} \end{cases}$$

Suppose D is a cube of dimension R .

- (a) Find the joint density (X_1, X_2, X_3) which is uniform on D .
- (b) Find the marginal density of X_1, X_2, X_3 .
- (c) Find the joint density of $(X_1, X_2), (X_1, X_3), (X_3, X_2)$.

Ex. 8.1.4. Let X_1, X_2, \dots, X_n be i.i.d. random variables having a common distribution function $F : \mathbb{R} \rightarrow [0, 1]$ and probability density function $f : \mathbb{R} \rightarrow \mathbb{R}$. Let $X_{(1)} < X_{(2)} < \dots < X_{(n)}$ be the corresponding order statistic. Show that for $1 \leq i < j \leq n$, $(X_{(i)}, X_{(j)})$ has a joint density function given by

$$f_{X_{(i)}, X_{(j)}}(x, y) = \frac{n!}{(i-1)!(j-1-i)!(n-j)!} f(x)f(y)[F(x)]^{i-1}[F(y)-F(x)]^{j-1-i}[1-F(y)]^{n-j},$$

for $-\infty < x < y < \infty$.

Ex. 8.1.5. Let X_1, X_2, \dots, X_n be i.i.d. random variables having a common distribution $X \sim \text{Uniform}(0, 1)$. Let $X_{(1)} < X_{(2)} < \dots < X_{(n)}$ be the corresponding order statistic. Show that $\frac{X_{(1)}}{X_{(n)}}$ and $X_{(n)}$ are independent random variables.

Ex. 8.1.6. Let $\{U_i : i \geq 1\}$ be a sequence of i.i.d. uniform $(0, 1)$ random variables and Let $N \sim \text{Poisson } (\lambda)$. Find the distribution of $V = \min\{U_1, U_2, \dots, U_{N+1}\}$.

Ex. 8.1.7. Let $-\infty < a < b < \infty$. Let X_1, X_2, \dots, X_n i.i.d $X \sim \text{Uniform}(a, b)$. Find the probability density function of $M = \frac{X_{(1)} + X_{(n)}}{2}$.

Ex. 8.1.8. Let X_1, X_2 be two independent standard normal random variables. Find the distribution of $Z = X_{(1)}^2$.

Ex. 8.1.9. Let X_1, X_2, \dots, X_n be i.i.d. Uniform $(0, 1)$ random variables.

- (a) Find the conditional distribution of $X_{(n)} | X_{(1)} = x$ for some $0 < x < 1$.
- (b) Find $E[X_{(n)} | X_{(1)} = x]$ and $\text{Var}[X_{(n)} | X_{(1)} = x]$.

Ex. 8.1.10. Suppose X is a symmetric continuous random variable. Let Y be a continuous random variable such that $P(Y > 0) = 1$. Show that $\frac{X}{Y}$ is symmetric.

Ex. 8.1.11. Verify (8.1.3).

Ex. 8.1.12. Suppose X_1, X_2, \dots are i.i.d. Cauchy $(0, 1)$ random variables.

- (a) Fix $z \in \mathbb{R}$. Find a, b, c, d such that

$$\frac{1}{1+x^2} \frac{1}{1+(z-x)^2} = \frac{ax+b}{1+x^2} + \frac{cx+d}{1+(z-x)^2},$$

for all $x \in \mathbb{R}$.

- (b) Show that $X_1 + X_2 \sim \text{Cauchy}(0, 2)$.
- (c) Use induction to show that $X_1 + X_2 + \dots + X_n \sim \text{Cauchy}(0, n)$.
- (d) Use Lemma 5.3.2 to show that $\bar{X}_n \sim \text{Cauchy}(0, 1)$.

8.2 WEAK LAW OF LARGE NUMBERS

Let $n \geq 1$, X_1, X_2, \dots, X_n be an i.i.d. random sample from a population whose distribution is given by a random variable X which has mean μ . In Chapter 7 we considered the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

and showed in Theorem 7.2.2 that $E[\bar{X}] = \mu$. We also discussed that \bar{X} could be considered as an estimate for μ . The below result makes this precise and is referred to as the weak law of large numbers.

In the statement and proof of the below Theorem we shall denote \bar{X} by \bar{X}_n to emphasise the dependence on n .

THEOREM 8.2.1. (Weak Law of Large Numbers) *Let X_1, X_2, \dots be a sequence of i.i.d. random variables. Assume that X_1 has finite mean μ and finite variance σ^2 . Then for any $\epsilon > 0$*

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \epsilon) = 0, \quad (8.2.1)$$

Proof- Let $\epsilon > 0$ be given. We note that

$$E(\bar{X}_n) = E\left(\frac{\sum_{i=1}^n X_i}{n}\right) = \sum_{i=1}^n \frac{1}{n} E(X_i) = \frac{n\mu}{n} = \mu.$$

Using Theorem 4.2.4, Theorem 4.2.6 and Exercise 6.2.17 we have

$$\begin{aligned} \text{Var}[\bar{X}_n] &= \text{Var}\left[\frac{\sum_{i=1}^n X_i}{n}\right] \\ &= \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n X_i\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] \\ &= \frac{\sigma^2}{n} \end{aligned}$$

So we have shown that the random variable \bar{X}_n has finite expectation variance. By Chebychev's inequality,(apply Theorem 6.1.13 (a) with $k = \frac{\epsilon}{\sigma}$), we have

$$P(|\bar{X}_n - \mu| > \epsilon) \leq \frac{\sigma^2}{n\epsilon^2}.$$

Therefore as $0 \leq P(|\bar{X}_n - \mu| > \epsilon)$ for all $n \geq 1$ and $\frac{\sigma^2}{n\epsilon^2} \rightarrow 0$ as $n \rightarrow \infty$, by standard results in Real Analysis we conclude that

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \epsilon) = 0.$$

■

REMARK 8.2.2. The convergence of sample mean to μ actually happens with Probability one. That is, suppose we denote the event $A = \{\lim_{n \rightarrow \infty} \bar{X}_n = \mu\}$, then $P(A) = 1$. The result is referred to as the Strong Law of large numbers. We prove it in Appendix C (see Theorem C.0.1).

Theorem 8.2.1 states that, for any $\epsilon > 0$, the $P(|\bar{X}_n - \mu| > \epsilon)$, goes to zero as $n \rightarrow \infty$. This mode of convergence of the sample to the true mean is called “convergence in probability”. We define it precisely below.

DEFINITION 8.2.3. A sequence X_1, X_2, \dots is said to converge in probability to a random variable X if for any $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0, \quad (8.2.2)$$

The following notation

$$X_n \xrightarrow{P} X$$

is typically used to convey that the sequence X_1, X_2, \dots converges in probability to X .

EXAMPLE 8.2.4. Let X_1, X_2, \dots, X_n be i.i.d random variables that are uniformly distributed over the interval $(0, 1)$. We already know by the law of large numbers that \bar{X} converges to $E(X_1) = \frac{1}{2}$ in probability. Often we are interested in other functionals of the sample and their convergence properties. We illustrate one such example below.

Consider the n -th order statistic $X_{(n)} = \max\{X_1, \dots, X_n\}$. For any $0 < \epsilon < 1$,

$$\begin{aligned} P(|X_{(n)} - 1| \geq \epsilon) &= P(X_{(n)} \leq 1 - \epsilon) + P(X_{(n)} \geq 1 + \epsilon) \\ &= P(X_{(n)} \leq 1 - \epsilon) + 0 \\ &= P\left(\cap_{i=1}^n (X_i \leq 1 - \epsilon)\right) \\ &= (1 - \epsilon)^n. \end{aligned}$$

and for $\epsilon > 1$,

$$P(|X_{(n)} - 1| \geq \epsilon) = P(X_{(n)} \leq 1 - \epsilon) + P(X_{(n)} \geq 1 + \epsilon) = 0.$$

For $0 < \epsilon < 1$,

$$\lim_{n \rightarrow \infty} (1 - \epsilon)^n = 0.$$

So we have shown that $X_{(n)}$ converges in probability to 1 as $n \rightarrow \infty$. ■

Another application of the weak law of large numbers is to sample proportion discussed in Section 7.2.3.

EXAMPLE 8.2.5. Suppose we are interested in an event A and want to estimate $p = P(X \in A)$. We consider a sample X_1, X_2, \dots, X_n which is i.i.d. X . We define a sequence of random variables $\{Y_n\}_{n \geq 1}$ by

$$Y_n = \begin{cases} 1 & \text{if } X_n \in A \\ 0 & \text{if } X_n \notin A \end{cases}$$

Clearly Y_n are independent (as the X_n are) and further they are identically distributed as $P(Y_n = 1) = P(X_n \in A) = p$. In particular $\{Y_n\}$ are an i.i.d. Bernoulli (p) sequence of random variables. We readily observe (as done in Chapter 7) that

$$\bar{Y}_n = \frac{\sum_{i=1}^n Y_i}{n} = \frac{\#\{X_i \in A\}}{n} = \hat{p}.$$

Hence the Weak law of large numbers (applied to the sequence Y_n) will imply that sample proportion converges to the true proportion p in probability. Consequently, as discussed earlier, this provides legitimacy to the relationship between Probability and relative frequency. ■

EXERCISES

Ex. 8.2.1. Let X, X_1, X_2, \dots, X_n be i.i.d random variables that are uniformly distributed over the interval $(0, 1)$. Consider the first order statistic $X_{(1)} = \max\{X_1, \dots, X_n\}$. Show that $X_{(1)}$ converges to 0 in probability.

Ex. 8.2.2. Let $X_1, X_2, \dots, X_n, \dots$ be i.i.d. random variables with finite mean and variance. Define

$$Y_n = \frac{2}{n(n+1)} \sum_{i=1}^n iX_i.$$

Show that $Y_n \xrightarrow{p} E(X_1)$ as $n \rightarrow \infty$.

8.3 CONVERGENCE IN DISTRIBUTION

When discussing a collection of random variables it makes sense to think of them as a sequence of objects, and as with any sequence in calculus we may ask whether the sequence converges in any way. We have already seen “convergence in probability” in the previous section. Here we be interested in what is known as “convergence in distribution”. This type of convergence plays a major role in the understand the limiting distribution of the sample mean (See Central Limit Theorem, Theorem 8.4.1).

DEFINITION 8.3.1. A sequence X_1, X_2, \dots is said to converge in distribution to a random variable X if $F_{X_n}(x)$ converges to $F_X(x)$ at every point x for which F_X is continuous. The following notation

$$X_n \xrightarrow{d} X$$

is typically used to convey that the sequence X_1, X_2, \dots converges in distribution to X .

EXAMPLE 8.3.2. Let $X_n \sim \text{Uniform}(0, \frac{1}{n})$ so that the distribution function is

$$F_{X_n}(x) = \begin{cases} 0 & \text{if } 0 \leq x \\ nx & \text{if } 0 < x < \frac{1}{n} \\ 1 & \text{if } x \geq \frac{1}{n} \end{cases}$$

and it is then easy to see that $F_{X_n}(x)$ converges to

$$F(x) = \begin{cases} 0 & \text{if } 0 \leq x \\ 1 & \text{if } x > 0 \end{cases}$$

If X is the constant random variable for which $P(X = 0) = 1$, then X has distribution function

$$F_X(x) = \begin{cases} 0 & \text{if } 0 < x \\ 1 & \text{if } x \geq 0 \end{cases}$$

It is not true that $F_X(x) = F(x)$, but the two are equal at points where they are continuous. Therefore the sequence X_1, X_2, \dots converges in distribution to the constant random variable 0. ■

Note that this form of convergence does not generally guarantee that probabilities associated with X can be derived as limits of probabilities associated with X_n . For instance, in the example above $P(X_n = 0) = 0$ for all n while $P(X = 0) = 1$. However, with a few additional assumptions a stronger claim may be made.

THEOREM 8.3.3. *Let f_{X_1}, f_{X_2}, \dots be the respective densities of continuous random variables X_1, X_2, \dots . Suppose they converge in distribution to a continuous random variable X with density f_X . Then for every interval A we have $P(X_n \in A) \rightarrow P(X \in A)$.*

Proof - Since X is a continuous random variable $F_X(x)$ is the integral of a density, and thus a continuous function. Therefore convergence in distribution guarantees that $F_{X_n}(x)$ converges to $F_X(x)$ everywhere. Let $A = (a, b)$ (and note that whether or not endpoints are included does not matter since all random variables are taken to be continuous). Then

$$\begin{aligned} P(X_n \in A) &= \int_a^b f_{X_n}(x) dx \\ &= F_{X_n}(b) - F_{X_n}(a) \\ &\rightarrow F_X(b) - F_X(a) \\ &= \int_a^b f_X(x) dx = P(X \in A). \end{aligned}$$
■

The second theorem about moment generating functions that we will state, but leave unproven, is the following:

THEOREM 8.3.4. (M.G.F. Convergence Theorem) *If X_1, X_2, \dots are a sequence of random variables whose moment generating functions $M_n(t)$ exist in an interval containing zero, and if $M_n(t) \rightarrow M(t)$ on that interval where $M(t)$ is the moment generating function of a random variable X , then X_n converges to X in distribution.*

To illustrate the use of this fact, consider an alternate proof of the limiting relationship between binomial and Poisson random variables (See Theorem 2.2.2).

EXAMPLE 8.3.5. Let $X \sim \text{Poisson}(\lambda)$ and let $X_n \sim \text{Binomial}(n, \frac{\lambda}{n})$. Then X_n converges in distribution to X .

The moment generating function of a binomial variable was already computed in Example 6.3.7. Therefore,

$$\begin{aligned} M_{X_n}(t) &= \left(\frac{\lambda}{n}e^t + (1 - \frac{\lambda}{n})\right)^n \\ &= \left(1 + \frac{\lambda(e^t - 1)}{n}\right)^n \end{aligned}$$

Using Exercise 8.4.4, we see that

$$M_{X_n}(t) \rightarrow e^{\lambda(e^t - 1)}.$$

On the other hand, the moment generating function of X is

$$\begin{aligned} M_X(t) &= E[e^{tX}] \\ &= \sum_{j=0}^{\infty} e^{tj} P(X = j) \\ &= \sum_{j=0}^{\infty} e^{tj} \frac{\lambda^j e^{-\lambda}}{j!} \\ &= e^{\lambda e^t} \cdot e^{-\lambda} \cdot \sum_{j=0}^{\infty} \frac{(\lambda e^t)^j e^{-\lambda e^t}}{j!} \\ &= e^{\lambda(e^t - 1)} \end{aligned}$$

where the series equals 1 since it is simply the sum of the probabilities of a Poisson(λe^t) random variable.

Since $M_{X_n}(t) \rightarrow M_X(t)$, by the m.g.f. convergence theorem (Theorem 8.3.4), X_n converges in distribution to X . That is, Binomial(n, p) random variables converge in distribution to a Poisson(λ) distribution when $p = \frac{\lambda}{n}$ and $n \rightarrow \infty$. ■

EXERCISES

Ex. 8.3.1. Suppose a sequence X_n , $n \geq 1$ of random variables converges to a random variable X in probability then show that X_n converges in distribution to X . That is show that

$$F_{X_n}(x) \rightarrow F_X(x) \text{ as } n \rightarrow \infty,$$

for all continuity points of $F_X : \mathbb{R} \rightarrow [0, 1]$ with F_{X_n}, F_X being the distribution functions of X_n and X respectively.

Ex. 8.3.2. Let X_n have the t -distribution with n degrees of freedom. Show that $X_n \xrightarrow{d} X$ where X is standard Normal distribution.

Ex. 8.3.3. Let $X_n \xrightarrow{d} X$. Show that $X_n^2 \xrightarrow{d} X^2$.

8.4 CENTRAL LIMIT THEOREM

Let $n \geq 1$, X_1, X_2, \dots, X_n be an i.i.d. random sample from a population with mean μ and variance σ^2 . Consider the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

As observed in Theorem 7.2.2, $E(\bar{X}) = \mu$ and $SD(\bar{X}) = \frac{\sigma}{\sqrt{n}}$. As discussed before, we might view this information as \bar{X} being typically close to μ up to an error of $\frac{\sigma}{\sqrt{n}}$ with high probability. As $n \rightarrow \infty$, $\frac{\sigma}{\sqrt{n}} \rightarrow 0$ and this indicates that \bar{X} approaches μ . We have already verified that \bar{X} converges in probability to μ

courtesy of the weak law of large numbers (in fact it converges with probability 1 by the strong law of large numbers).

To get a better understanding of the limiting distribution of \bar{X} we standardise it and consider,

$$Y_n = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}.$$

Finding the probabilities of events connected with Y_n for each n exactly may not be possible in all cases but one can find good approximate values. It turns out that for a large class of random variables the distribution of Y_n is close to that of the standard Normal random variable particularly for large n . This remarkable fact is referred to as the Central Limit Theorem and we prove it next.

As done earlier, in the statement and proof of the below Theorem we shall denote \bar{X} by \bar{X}_n to emphasise the dependence on n .

THEOREM 8.4.1. (Central Limit Theorem) *Let X_1, X_2, \dots be i.i.d. random variables with finite mean μ , finite variance σ^2 , and possessing common moment generating function $M_X()$. Then*

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} Z, \quad (8.4.1)$$

where $Z \sim \text{Normal}(0, 1)$.

Proof- Let $Y_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}$. We will verify that

$$\lim_{n \rightarrow \infty} M_{Y_n}(t) = e^{\frac{t^2}{2}}.$$

Now, using the definition of the moment generating function and some elementary algebra we have

$$\begin{aligned} M_{Y_n}(t) &= E[\exp(tY_n)] = E\left[\exp\left(t\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}\right)\right] \\ &= E\left[\exp\left(\frac{t}{\sigma}\sqrt{n}\left(\frac{\sum_{i=1}^n X_i}{n} - \mu\right)\right)\right] = E\left[\exp\left(\sum_{i=1}^n \frac{t}{\sigma\sqrt{n}}(X_i - \mu)\right)\right] \\ &= E\left[\prod_{i=1}^n \exp\left(\frac{t}{\sigma\sqrt{n}}(X_i - \mu)\right)\right]. \end{aligned} \quad (8.4.2)$$

As X_1, X_2, \dots, X_n are independent, from Theorem 8.1.2 we can conclude that

$$\exp\left(\frac{t}{\sigma\sqrt{n}}(X_1 - \mu)\right), \exp\left(\frac{t}{\sigma\sqrt{n}}(X_2 - \mu)\right), \dots, \exp\left(\frac{t}{\sigma\sqrt{n}}(X_n - \mu)\right)$$

are also independent. From Exercise 7.2.2 and 7.2.3, they also have the same distribution. So from the calculation in (8.4.2) and using Exercise 6.3.4 inductively we have

$$\begin{aligned} M_{Y_n}(t) &= E\left[\prod_{i=1}^n \exp\left(\frac{t}{\sigma\sqrt{n}}(X_i - \mu)\right)\right] = \prod_{i=1}^n E\left[\exp\left(\frac{t}{\sigma\sqrt{n}}(X_i - \mu)\right)\right] \\ &\quad (\text{Using Theorem 6.3.9(a)}) \\ &= \left(E\left[\exp\left(\frac{t}{\sigma\sqrt{n}}(X_1 - \mu)\right)\right]\right)^n. \end{aligned} \quad (8.4.3)$$

Let $U = \frac{X_1 - \mu}{\sigma}$. As $E(U) = 0, E(U^2) = 1$ we have that $M'_U(0) = 0$ and $M''_U(0) = 1$. From Exercise 8.4.5, we have that for $t \in \mathbb{R}$

$$M_U(t) = 1 + \frac{t^2}{2} + g(t) \quad (8.4.4)$$

where $\lim_{s \rightarrow 0} \frac{g(s)}{s^2} = 0$. Therefore from (8.4.3) and (8.4.4) we have

$$M_{Y_n}(t) = [M_U(t)]^n = \left[1 + \frac{t^2}{2n} + g\left(\frac{t}{\sqrt{n}}\right)\right]^n = \left[1 + \frac{1}{n}\left(\frac{t^2}{2} + ng\left(\frac{t}{\sqrt{n}}\right)\right)\right]^n$$

Using the fact $\frac{t^2}{2} + ng\left(\frac{t}{\sqrt{n}}\right) \rightarrow \frac{t^2}{2}$ and Exercise 8.4.4 it follows that ,

$$\lim_{n \rightarrow \infty} M_{Y_n}(t) = e^{\frac{t^2}{2}}.$$

Theorem 8.3.4 will then imply the result. ■

REMARK 8.4.2. The existence of moment generating function is not essential for the Central Limit Theorem. (8.4.1) holds when X, X_1, X_2, \dots are i.i.d. random variables with finite mean μ and finite variance σ^2 . The proof is more complicated in this case.

Further we shall often use an equivalent formulation of (8.4.1). By definition of \bar{X} and elementary algebra we see that $Y_n = \frac{S_n - n\mu}{\sqrt{n}\sigma}$, where $S_n = \sum_{i=1}^n X_i$.

$$\frac{S_n - n\mu}{\sqrt{n}\sigma} \xrightarrow{d} Z, \quad (8.4.5)$$

where $S_n = \sum_{i=1}^n X_i$.

8.4.1 Normal Approximation and Continuity Correction

A typical application of the central limit theorem is to find approximate value of the probability of events related to S_n or \bar{X} . For instance, suppose we were interested in calculating for any $a, b \in \mathbb{R}$, $P(a < S_n \leq b)$ for large n . We would proceed in the following way. We know from (8.4.5) that

$$P\left(\frac{S_n - n\mu}{\sqrt{n}\sigma} \leq x\right) \rightarrow P(Z \leq x) \quad (8.4.6)$$

as $n \rightarrow \infty$ for all $x \in \mathbb{R}$.

$$\begin{aligned} P(a < S_n \leq b) &= P\left(\frac{a - n\mu}{\sqrt{n}\sigma} < \frac{S_n - n\mu}{\sqrt{n}\sigma} \leq \frac{b - n\mu}{\sqrt{n}\sigma}\right) \\ &= P\left(\frac{S_n - n\mu}{\sqrt{n}\sigma} \leq \frac{b - n\mu}{\sqrt{n}\sigma}\right) - P\left(\frac{S_n - n\mu}{\sqrt{n}\sigma} \leq \frac{a - n\mu}{\sqrt{n}\sigma}\right) \\ &\quad \text{from (8.4.6) for large enough } n \\ &\approx P\left(Z \leq \frac{b - n\mu}{\sqrt{n}\sigma}\right) - P\left(Z \leq \frac{a - n\mu}{\sqrt{n}\sigma}\right) \\ &= P\left(\frac{a - n\mu}{\sqrt{n}\sigma} < Z \leq \frac{b - n\mu}{\sqrt{n}\sigma}\right), \end{aligned}$$

where in the second last line we have used the notation \approx to indicate that the right hand side is an approximation. Therefore we would conclude that for large n ,

$$P(a < S_n \leq b) \approx P\left(\frac{a - n\mu}{\sqrt{n}\sigma} < Z \leq \frac{b - n\mu}{\sqrt{n}\sigma}\right). \quad (8.4.7)$$

We would then use the R function `pnorm()` or Normal Tables (See Table D.2) to compute the right hand side.

A similar computation would also yield

$$P(a < \bar{X} \leq b) \approx P\left(\frac{\sqrt{n}(a - \mu)}{\sigma} < Z \leq \frac{\sqrt{n}(b - \mu)}{\sigma}\right). \quad (8.4.8)$$

EXAMPLE 8.4.3. Let Y be a random variable distributed as $\text{Gamma}(100, 4)$. Suppose we were interested in finding

$$P(20 < Y \leq 30).$$

Suppose X_1, X_2, \dots, X_{100} are independent Exponential (4) random variables then Y and $S_{100} = \sum_{i=1}^{100} X_i$ have the same distribution. Therefore, applying the Central Limit Theorem with $\mu = E(X_1) = \frac{1}{4}$, $\sigma = \text{SD}(X_1) = \frac{1}{4}$, we have

$$\begin{aligned} P(20 < Y \leq 30) &= P(20 < S_{100} \leq 30) \\ &\quad \text{by (8.4.7)} \\ &\approx P\left(\frac{20 - 100(0.25)}{\sqrt{100}(0.25)} < Z \leq \frac{30 - 100(0.25)}{\sqrt{100}(0.25)}\right) \\ &= P\left(\frac{-5}{2.5} < Z \leq \frac{5}{2.5}\right) \\ &= P(-2 < Z \leq 2) \\ &= P(Z \leq 2) - P(Z \leq -2) \\ &\quad \text{using symmetry of Normal distribution} \\ &= P(Z \leq 2) - (1 - P(Z \leq 2)) \\ &= 2P(Z \leq 2) - 1 \end{aligned}$$

Looking up Table D.2, we see that this value comes out to be approximately $2 \times 0.9772 - 1 = 0.9544$. A more precise answer is given by R as

```
> 2 * pnorm(2) - 1
[1] 0.9544997
```

Using R, we can also compare this with the exact probability that we are approximating.

```
> pgamma(30, 100, 4) - pgamma(20, 100, 4)
[1] 0.9550279
```

■

Continuity Correction: Suppose X_1, X_2, X_3, \dots are all integer valued random variables. Then $S_n = \sum_{i=1}^n X_i$ is also a integer random variable. Now,

$$P(S_n = k) = P(k - h < S_n \leq k + h)$$

for all natural numbers k and $0 < h < 1$. However it is easy to see that two distinct values of h will lead to two different answers if we use the Normal approximation provided by the Central Limit Theorem. One can also observe that this will increase with h . So it is customary to use $h = \frac{1}{2}$ while computing such probabilities using the Normal approximation. So when X_1, X_2, X_3, \dots are all integer valued random variables we use,

$$\begin{aligned} P(a < S_n \leq b) &= P(a - 0.5 < S_n \leq b + 0.5) \\ &\approx P\left(\frac{a + 0.5 - n\mu}{\sqrt{n}\sigma} < Z \leq \frac{b + 0.5 - n\mu}{\sqrt{n}\sigma}\right) \end{aligned} \tag{8.4.9}$$

whenever a, b are in the range of S_n . This convention is referred to as the “continuity correction”.

EXAMPLE 8.4.4. Two types of coin are produced at a factory: a fair coin and a biased one that comes up heads 55% of the time. Priya is the quality control scientist at the factory. She wants to design an experiment that will test whether a coin is fair or biased. In order to ascertain which type of coin she has, she prescribes the following experiment as a test:- *Toss the given coin 1000 times, if the coin comes up heads 525 or more times conclude that it is a biased coin. Otherwise conclude that it is fair.* Factory manager Ayesha is interested in the following question: What is the probability that Priya’s test shall reach a false conclusion for a fair coin ?

Let S_{1000} be the number of heads in 1000 tosses of a coin. As discussed in earlier chapters, we know that $S_{1000} = \sum_{i=1}^{1000} X_i$ where each X_i are i.i.d. Bernoulli random variables with parameter p .

If the coin is fair, then $p = 0.5$ and $E[X_1] = 0.5$, $Var[X_1] = 0.25$, and therefore $E[S_{1000}] = 500$ and $SD[S_{1000}] = \sqrt{250} = 15.8114$. We want to approximate

$$P(S_{1000} \geq 525) = 1 - P(S_{1000} \leq 524) = 1 - P(S_{1000} \leq 524.5)$$

Without the continuity correction, we would approximate this probability as

$$1 - P\left(Z \leq \frac{24}{15.8114}\right) = 1 - P(Z \leq 1.52)$$

which can be computed using Table D.2 as $1 - 0.9357 = 0.0643$, or using R as

```
> 1 - pnorm(24 / sqrt(250))
[1] 0.06452065
```

With the continuity correction, the approximate value would instead use $z = 24.5/15.8114 = 1.55$, giving $1 - 0.9394 = 0.0606$ using Table D.2 or

```
> 1 - pnorm(24.5 / sqrt(250))
[1] 0.06062886
```

in R. We can also compute the exact probability that we are trying to approximate, namely $P(S_{1000} \geq 525)$, in R as

```
> 1 - pbinom(524, 1000, 0.5)
[1] 0.06060713
```

As we can see, the continuity correction gives us a slightly better approximation. These calculations tell us that the probability of Priya's test reaching a false conclusion if the coin is fair is approximately 0.061. We shall examine the topic of Hypothesis testing, that Priya was trying to do, more in Chapter 9. ■

EXAMPLE 8.4.5. We return to the Birthday problem. Suppose a small town has 1095 students. What is the probability that five or more students were born on independence day? Assume that birthrates are constant throughout the year and that each year has 365 days.

The probability that any given student was born on independence day is $\frac{1}{365}$. So the exact probability that five or more students were born on independence day is

$$1 - \sum_{k=0}^4 \binom{1095}{k} \left(\frac{1}{365}\right)^k \left(\frac{364}{365}\right)^{1095-k}.$$

In Example 2.2.1 we have used the Poisson approximation with $\lambda = 4$ to estimate the above as

$$\begin{aligned} & 1 - \sum_{k=0}^4 \binom{1460}{k} \left(\frac{1}{365}\right)^k \left(\frac{364}{365}\right)^{1460-k} \\ & \approx 1 - \left[e^{-4} + 4e^{-4} + \frac{4^2}{2}e^{-4} + \frac{1}{6}4^3e^{-4} + \frac{1}{24}4^4e^{-4} \right] \\ & = 0.3711631 \end{aligned}$$

We can do another approximation using central limit theorem, which is typically called the normal approximation. For $1 \leq i \leq 1460$, define

$$X_i = \begin{cases} 1 & \text{if } i\text{th person's birthday is on independence day} \\ 0 & \text{otherwise} \end{cases}$$

Given the assumptions above on birthrates we know X_i are i.i.d random variables distributed as Bernoulli($\frac{1}{365}$). Note that $S_{1460} = \sum_{i=1}^{1460} X_i$ is the number of people born on independence day and we are interested in calculating

$$P(S_{1460} \geq 5).$$

Observe that $E(X_1) = \frac{1}{365}$, $\text{Var}(X_1) = \frac{1}{365}(1 - \frac{1}{365}) = \frac{364}{365^2}$. By the central limit theorem, we know that

$$\begin{aligned} P(S_{1460} \geq 5) &= 1 - P(S_{1460} \leq 4) = 1 - P(S_{1460} \leq 4.5) \\ &\approx 1 - P(Z \leq \frac{4.5 - (1460)(\frac{1}{365})}{\sqrt{(1460)(\frac{364}{365^2})}}) \\ &= 1 - P(Z \leq \frac{0.5}{1.9973}) \\ &= 0.401 \end{aligned}$$

Recall from the calculations done in Example 2.2.1 that the exact answer for this problem is 0.3711629. So in this example, the Poisson approximation seems to work better than the Normal approximation. This is due to the fact that more asymmetry in the underlying Bernoulli distribution worsens the normal approximation, just as it improves the Poisson approximation as we saw in Figure 2.2. ■

EXERCISES

Ex. 8.4.1. Suppose S_n is binomially distributed with parameters $n = 200$ and $p = 0.3$. Use the central limit theorem to find an approximation for $P(99 \leq S_n \leq 101)$.

Ex. 8.4.2. Toss a fair coin 400 times. Use the central limit theorem to

- (a) find an approximation for the probability of at most 190 heads.
- (b) find an approximation for the probability of at least 70 heads.
- (c) find an approximation for the probability of at least 120 heads.
- (d) find an approximation for the probability that the number of heads is between 140 and least 160.

Ex. 8.4.3. Suppose that the weight of open packets of daal in a home is uniformly distributed from 200 to 600 gms. In random survey of 64 homes, find the (approximate) probability that the total weight of open boxes is less than 25 kgs.

Ex. 8.4.4. Let $\{a_n\}_{n \geq 1}$ be a sequence of real numbers such that $a_n \rightarrow a$ as $n \rightarrow \infty$. Then

$$\lim_{n \rightarrow \infty} \left(1 + \frac{a_n}{n}\right)^n = e^a.$$

Ex. 8.4.5. Suppose U is a random variable (discrete or continuous) and $M_U(t) = E(e^{tU})$ exists for all t . Then show that

$$M_U(t) = 1 + tM'_U(0) + \frac{t^2}{2}M''_U(0) + g(t)$$

where $\lim_{t \rightarrow 0} \frac{g(t)}{t^2} = 0$.

Ex. 8.4.6. Let $\{X_n\}_{n \geq 1}$ be a sequence of i.i.d. random variables with $X_1 \sim \text{Exponential}(1)$. Find

$$\lim_{n \rightarrow \infty} P\left(\frac{n}{2} - \frac{\sqrt{n}}{2\sqrt{3}} \leq \sum_{i=1}^n [1 - \exp(-X_i)] \leq \frac{n}{2} + \frac{\sqrt{n}}{2\sqrt{3}}\right).$$

Ex. 8.4.7. Let $a_n = \sum_{k=0}^n \frac{n^k}{k!} e^{-n}$, $n \geq 1$. Using the Central Limit Theorem evaluate $\lim_{n \rightarrow \infty} a_n$.

Ex. 8.4.8. How often should you toss a coin:

- (a) to be at least 90 % sure that your estimate of the P(head) is within 0.1 of its true value ?
- (b) to be at least 90 % sure that your estimate of the P(head) is within 0.01 of its true value ?

Ex. 8.4.9. To forecast the outcome of the election in which two parties are contesting, an internet poll via Facebook is conducted. How many people should be surveyed to be at least 95% sure that the estimated proportion is within 0.05 of the true value ?

Ex. 8.4.10. A medical study is conducted to estimate the proportion of people suffering from April allergies in Bangalore. How many people should be surveyed to be at least 99% sure that the estimate is within 0.02 of the true value ?

ESTIMATION AND HYPOTHESIS TESTING

In Chapter 7 we introduce the question of how an i.i.d. sample X_1, X_2, \dots, X_n from an unknown distribution may be used to estimate aspects of that distribution. In Chapter 8 we saw how the sample statistics behave asymptotically. In this chapter we look at some specific examples where various parameters of the distribution such as μ and σ are unknown, and the sample is used to estimate these parameters.

For example, suppose there is a coin which we assume has a probability p of showing heads each time it is flipped. To gather information about p the coin is flipped 100 times. The results of these flips are viewed as i.i.d. random variables X_1, X_2, \dots, X_{100} with a $\text{Bernoulli}(p)$ distribution. Suppose $\sum_{n=1}^{100} X_n = 60$, meaning 60 of the 100 flips showed heads. How might we use this to infer something about the value of p ?

The first two topics we will consider are the “method of moments” and the “maximum likelihood estimate”. Both of these are direct forms of estimation in the sense that they produce a single-value estimate for p . A benefit of such methods is that they produce a single prediction, but a downside is that the prediction they make is most likely not exactly correct. These methods amount to a statement like “Since 60 of the 100 flips came up heads, we predict that the coin should come up heads 60% of the time in the long run”. In some sense the 60% prediction may be the most reasonable one given what was observed in the 100 flips, but it should also be recognised that 0.6 is unlikely to be the true value of p .

Another approach is that of the “confidence interval”. Using this method we abandon the hope of realising a specific estimate and instead produce a range of values in which we expect to find the unknown parameter. This yields a statement such as, “With 90% confidence the actual probability the coin will show heads is between 0.52 and 0.68”. While this approach does not give a single-valued estimate, it has the benefit that the result is more likely to be true.

Yet another approach is the idea of a “hypothesis test”. In this case we make a conjecture about the value of the parameter and make a computation to test the credibility of the conjecture. The result may be a statement such as, “If the coin had a 50% chance of showing heads, then observing 60 heads or more in 100 flips should occur less than 3% of the time. This is a rare enough result, it suggests that the 50% hypothesis of showing heads is inaccurate”.

For all of these methods, we will assume that the sample X_1, X_2, \dots, X_n are i.i.d copies of a random variable X with a probability mass function or probability density function $f(x)$. For brevity, we shall often refer to the distribution X , by which we will mean the distribution of the random variable X . We shall further assume that $f(x)$ depends on one or more unknown parameters p_1, p_2, \dots, p_d and emphasise this using the notation $f(x | p_1, p_2, \dots, p_d)$. We may abbreviate this $f(x | p)$ where p represents the vector of all parameters $(p_1, \dots, p_d) \in \mathcal{P} \subset \mathbb{R}^d$ for some $d \geq 1$. The set \mathcal{P} may be all of \mathbb{R}^d or some proper subset depending on the nature of the parameters.

9.1 NOTATIONS AND TERMINOLOGY FOR ESTIMATORS

DEFINITION 9.1.1. Let $X_1, X_2, X_3, \dots, X_n$ be an i.i.d. sample from the population with distribution X . Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$. Then $g(X_1, X_2, \dots, X_n)$ is defined as a “point estimator” from the sample and the value from a particular realisation is called an “estimate”.

In practice the function g is chosen keeping in mind the parameter of interest. We have seen the following in Chapter 7.

EXAMPLE 9.1.2. Let $E[X] = \mu$. Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by

$$g(x) = \frac{1}{n} \sum_{i=1}^n x_i.$$

Then $g(X_1, X_2, \dots, X_n)$ is the (now familiar) sample mean and it is an estimator for μ . We also saw that $E[g(X_1, X_2, \dots, X_n)] = \mu$ regardless of the true value of μ and we called such an estimator an unbiased estimator. ■

As noted in Chapter 7, we can view this as estimating the first moment of a distribution by the first moment of the empirical distribution based on a sample. A generalization of this method is known as the method of moments.

9.2 METHOD OF MOMENTS

Let X_1, X_2, \dots, X_n be a sample with distribution X . Assume that X is either has probability mass function or probability density function $f(x | p)$ depending on parameter(s) $p = (p_1, \dots, p_d)$. For $d \geq 1$. Let $m_k : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by

$$m_k(x) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

Notice that $m_k(X_1, X_2, \dots, X_n)$ is the k -th moment of the empirical distribution based on the sample X_1, X_2, \dots, X_n , which we will refer to simply as the k -th moment of the sample.

Let $\mu_k = E[X^k]$, the k -th moment of the distribution X . Since the distribution of X depends on (p_1, p_2, \dots, p_d) one can view μ_k as a function of p , which we can make explicit by the notation $\mu_k \equiv \mu_k(p_1, p_2, \dots, p_d)$. The method of moments estimator for (p_1, p_2, \dots, p_d) is obtained by equating the first d moments of the sample to the corresponding moments of the distribution. Specifically, it requires solving the d equations in d unknowns given by

$$\mu_k(p_1, p_2, \dots, p_d) = m_k(X_1, X_2, \dots, X_n), \quad k = 1, 2, \dots, d.$$

for p_1, p_2, \dots, p_d . There is no guarantee in general that these equations have a unique solution or that it can be computed, but in practice it is usually possible to do so. The solution will be denoted by $\hat{p}_1, \hat{p}_2, \dots, \hat{p}_d$ which will be written in terms of the realised values for $m_k, k = 1, \dots, d$. We will now explore this method for two examples.

EXAMPLE 9.2.1. Suppose X_1, X_2, \dots, X_{10} is an i.i.d. sample with distribution Binomial(N, p) where neither N nor p is known. Suppose the empirical realisation of these variables is 8, 7, 6, 11, 8, 5, 3, 7, 6, 9. One can check that the average of these values is $m_1 = 7$ while the average of their squares is $m_2 = 53.4$. Since $X \sim \text{Binomial}(N, p)$ the probability mass function is given by ???. We have previously shown that

$$E[X] = Np \text{ and } E[X^2] = \text{Var}[X] + E[X]^2 = Np(1-p) + N^2p^2.$$

Thus, the method of moments estimator for (N, p) is obtained by solving

$$7 = m_1 = \hat{N}\hat{p} \text{ and } 53.4 = m_2 = \hat{N}\hat{p}(1-\hat{p}) + \hat{N}^2\hat{p}^2.$$

Using elementary algebra we see that

$$\begin{aligned} \hat{N} &= \frac{m_1^2}{m_1 - (m_2 - m_1^2)} \approx 19 \\ \hat{p} &= \frac{m_1 - (m_2 - m_1^2)}{m_1} \approx 0.371. \end{aligned}$$

The method of moments estimates that the distribution from which the sample came is Binomial(19, 0.371). As we noted at the beginning of the chapter, we may wish to restrict the parameters based on the context of the problem. Since the N value is surely some integer, the estimate of \hat{N} was rounded to the nearest meaningful value in this case. ■

EXAMPLE 9.2.2. Suppose our quantity of interest X has a Normal (μ, σ^2) distribution. Therefore our probability density function is given by

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

Let $X_1, X_2, X_3, \dots, X_n$ be an i.i.d. sample from the population. We have shown that

$$E[X] = \mu \text{ and } E[X^2] = \text{Var}[X] + E[X]^2 = \mu^2 + \sigma^2.$$

The method of moments estimator for μ, σ is found by solving

$$m_1 = \mu \text{ and } m_2 = \mu^2 + \sigma^2.$$

from which

$$\begin{aligned} \hat{\mu} &= m_1 = \bar{X} \text{ and} \\ \hat{\sigma} &= \sqrt{m_2 - m_1^2} = \sqrt{\frac{n-1}{n}} S. \end{aligned}$$

■

It happens that the method of moment estimators may not be very reliable. For instance in the first example the estimate for p could be negative, occurring when the sample mean is smaller than the sample variance. Such defects can be somewhat rectified using moment matching and other techniques (see [CasBer90]).

9.3 MAXIMUM LIKELIHOOD ESTIMATE

Let $n \geq 1$, $p = (p_1, p_2, \dots, p_d) \in \mathbb{R}^d$ and X_1, X_2, \dots, X_n be a sample from the population described by X . Assume that X either has probability mass function or probability density function denoted by $f(x | p)$ depending on parameter(s) $p \in \mathcal{P} \subset \mathbb{R}^d$.

DEFINITION 9.3.1. The likelihood function for the sample (X_1, X_2, \dots, X_n) is the function $L : \mathcal{P} \times \mathbb{R}^n \rightarrow \mathbb{R}$ given by

$$L(p; X_1, \dots, X_n) = \prod_{i=1}^n f(X_i | p).$$

For a given (X_1, X_2, \dots, X_n) , suppose $\hat{p} \equiv \hat{p}(X_1, X_2, \dots, X_n)$ is the point at which $L(p; X_1, \dots, X_n)$ attains its maximum as a function p . Then \hat{p} is called the maximum likelihood estimator of p (or abbreviated as MLE of p) given the sample (X_1, X_2, \dots, X_n) .

One observes readily that the likelihood function is the joint density or joint mass function of (X_1, X_2, \dots, X_n) . The MLE \hat{p} obtained is the most likely value of the parameter p , given that it is the value at which f is maximised for the given realisation (X_1, X_2, \dots, X_n) .

EXAMPLE 9.3.2. Let $p \in \mathbb{R}$ and (X_1, X_2, \dots, X_n) be from a population distributed as Normal with mean p and variance 1. Then the likelihood function is given by

$$L(p; X_1, \dots, X_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(X_i-p)^2}{2}} = \frac{1}{(\sqrt{2\pi})^n} e^{-\sum_{i=1}^n \frac{(X_i-p)^2}{2}}.$$

To find the MLE, treating the given the realisation X_1, X_2, \dots, X_n as fixed, one needs to maximise L as a function of p . This is equivalent to finding the minimum of $g : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$g(p) = \sum_{i=1}^n (X_i - p)^2.$$

Method 1: Since $g(p) = \sum_{i=1}^n (X_i - \bar{X})^2 + (\bar{X} - p)^2$ (See Exercise 9.3.1) and first term is always non-negative, the minimum of g will occur when

$$p = \bar{X}.$$

Method 2: The second method is to find the MLE using differential calculus. As g is a quadratic in p , it is differentiable at all p and

$$g'(p) = -2 \sum_{i=1}^n (X_i - p).$$

As the coefficient of p^2 in g is n (which is positive), and g is quadratic, the minimum will occur in the interior when $g'(p) = 0$. This occurs when p is equal to $\frac{1}{n} \sum_{i=1}^n X_i$. So the MLE of p is given by

$$\hat{p} = \bar{X}.$$

■

EXAMPLE 9.3.3. Let $p \in (0, 1)$ and (X_1, X_2, \dots, X_n) be from a population distributed as Bernoulli (p). Now the probability mass function f can be written as

$$f(x | p) = \begin{cases} p & \text{if } x = 1 \\ 1-p & \text{if } x = 0 \\ 0 & \text{otherwise.} \end{cases} = \begin{cases} p^x (1-p)^{1-x} & \text{if } x \in \{0, 1\} \\ 0 & \text{otherwise.} \end{cases}$$

Then the likelihood function is given by

$$L(p; X_1, \dots, X_n) = \prod_{i=1}^n p^{X_i} (1-p)^{1-X_i} = p^{\sum_{i=1}^n X_i} (1-p)^{n-\sum_{i=1}^n X_i}.$$

To find the MLE, treating the given realisation X_1, X_2, \dots, X_n as fixed, one needs to maximise L as a function of p . One needs to use calculus to find the MLE but differentiating L is cumbersome. So we will look at the logarithm of L (called the log likelihood function).

$$\begin{aligned} T(p; X_1, \dots, X_n) &= \ln L(p; X_1, \dots, X_n) \\ &= \begin{cases} \ln \left(\frac{p}{1-p} \right) a + n \ln(1-p) & \text{if } \sum_{i=1}^n X_i = a, 0 < a < n \\ n \ln(1-p) & \text{if } \sum_{i=1}^n X_i = 0 \\ n \ln(p) & \text{if } \sum_{i=1}^n X_i = n \end{cases} \end{aligned}$$

Therefore, in the first case, differentiating and setting it to zero $p = \frac{a}{n}$. This in fact can be verified to be the global maximum. In the second case T is a decreasing function of p and maximum occurs at $p = 0$. In the final case T is an increasing function of p and maximum occurs at $p = 1$. Therefore we can conclude that

$$\hat{p} = \frac{\sum_{i=1}^n X_i}{n}.$$

■

As a final example, let us revisit Example 9.2.1, where we considered a Binomial distribution with both parameters unknown.

EXAMPLE 9.3.4. Suppose X_1, X_2, \dots, X_n is an i.i.d. sample with distribution $\text{Binomial}(N, p)$ where neither N nor p is known. The likelihood function is given by

$$L(N, p; X_1, \dots, X_n) = \prod_{i=1}^k \binom{N}{x_i} p^{X_i} (1-p)^{N-X_i}$$

...

EXERCISES

Ex. 9.3.1. Show that for any real numbers p, x_1, x_2, \dots, x_n

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

9.4 CONFIDENCE INTERVALS

Let X_1, X_2, \dots, X_n be an i.i.d. sample from a distribution X with unknown mean μ . The sample mean \bar{X} is an unbiased estimator for μ , but the empirical value of the sample mean will likely differ from μ by some unknown amount. Suppose we want to produce an interval, centered around \bar{X} which we can be fairly certain contains the true average μ . This is known as a “confidence interval” and we explore how to produce such a thing in two different settings below.

9.4.1 *Confidence Intervals when the standard deviation σ is known*

Let X have a probability mass function or probability density function $f(x | \mu)$ where the distribution X has an unknown expected value μ , but a known standard deviation σ . Let X_1, X_2, \dots, X_n be an i.i.d. sample from the distribution X . Let $\beta \in (0, 1)$ denote a “confidence level”. We want to find an interval width a such that

$$P(|\bar{X} - \mu| < a) = \beta.$$

That is, the sample mean \bar{X} will have a probability β of differing from the true mean μ by no more than the quantity a .

Let

$$Z = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}.$$

Then $E[Z] = 0$ and $Var[Z] = 1$. Observe,

$$\begin{aligned} \beta &= P(|\bar{X} - \mu| < a) = P\left(\left|\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}\right| < \frac{\sqrt{n}a}{\sigma}\right) \\ &= P\left(|Z| < \frac{\sqrt{n}a}{\sigma}\right) \\ &= P\left(-\frac{\sqrt{n}a}{\sigma} < Z < \frac{\sqrt{n}a}{\sigma}\right) \end{aligned}$$

If X has a normal distribution, then $Z \sim Normal(0, 1)$ by Example 6.3.12. If the distribution of X is unknown, but n is large, then by the Central Limit Theorem (Theorem 8.4.1), Z should still be roughly normal, so we can make a valid assumption that $Z \sim Normal(0, 1)$. When n , σ , and β are known, one can use the Normal tables (Table D.2) to find the unknown interval width a . The interval $(\bar{X} - a, \bar{X} + a)$ is then known as “a β confidence interval for μ ”. The interpretation being that the random sample of size n from the distribution X should produce confidence intervals that include the correct value of μ 95% with probability β .

EXAMPLE 9.4.1. Suppose X has a normal distribution with known standard deviation $\sigma = 3.0$ and an unknown mean μ . A sample X_1, X_2, \dots, X_{16} of i.i.d. random variables is taken from distribution X . The sample average of these 16 values comes out to be $\bar{X} = 10.2$. What would be a 95% confidence interval for the actual mean μ ?

In this case $\beta = 0.95$ so we must find the value of a for which

$$P(|\bar{X} - \mu| < a) = 0.95.$$

From the computation above, this is equivalent to the equation

$$P\left(-\frac{4a}{3} < Z < \frac{4a}{3}\right) = 0.95$$

where $Z \sim \text{Normal}(0, 1)$. Using the normal table, this is equivalent to $\frac{4a}{3} \approx 1.96$, and so $a \approx 1.47$. In other words, a 95% confidence interval for the actual mean of the distribution X is $(8.73, 11.67)$.

It should be noted that the only random variable in the expression $P(|\bar{X} - \mu| < a) = 0.95$ is the \bar{X} variable. The interpretation is that random samples of size $n = 16$ from the distribution X should produce confidence intervals that include the correct value of μ , 95% of the time. ■

9.4.2 Confidence Intervals when the standard deviation σ is unknown

In most realistic situations the standard deviation σ would be unknown and would have to be estimated from the sample standard deviation S . In this case a confidence interval may still be produced, but an approximation via a normal distribution is insufficient.

Let X have normal distribution with density function $f(x | \mu, \sigma)$ where μ and σ are unknown. Let X_1, X_2, \dots, X_n be an i.i.d. sample from the distribution X . Let \bar{X} be the sample mean and the sample variance be denoted by S^2 . Let $\beta \in (0, 1)$ denote a confidence level. As before, we want to find an interval width a such that $P(|\bar{X} - \mu| < a) = \beta$.

Let

$$T = \frac{\sqrt{n}(\bar{X} - \mu)}{S}.$$

From Corollary 8.1.10, $T \sim t_{n-1}$. In a similar fashion as the previous example,

$$\begin{aligned} \beta &= P(|\bar{X} - \mu| < a) = P\left(\left|\frac{\sqrt{n}(\bar{X} - \mu)}{S}\right| < a\right) \\ &= P\left(|T| < \frac{\sqrt{n}a}{S}\right) \\ &= P\left(-\frac{\sqrt{n}a}{S} < T < \frac{\sqrt{n}a}{S}\right) \end{aligned}$$

where $T \sim t_{n-1}$. Since n , S , and β are known, this equation can be solved to find the unknown interval width a using the t-distribution.

EXAMPLE 9.4.2. Suppose X has a normal distribution with unknown mean μ and unknown standard deviation σ . A sample X_1, X_2, \dots, X_{16} of i.i.d. random variables with distribution X is taken. The sample average of these 16 values comes out to be $\bar{X} = 10.2$ and the sample standard deviation is $S = 3.0$. What would be a 95% confidence interval for the actual mean μ ?

In this case $\beta = 0.95$ so we must find the value of a for $P(|\bar{X} - \mu| < a) = 0.95$. This is equivalent to the equation

$$P\left(-\frac{4a}{3} < T < \frac{4a}{3}\right) = 0.95$$

Using the t-distribution, this is equivalent to $\frac{4a}{3} \approx 2.13$, and so $a \approx 1.60$. In other words, a 95% confidence interval for the actual mean of the distribution X is $(8.6, 11.80)$.

It is useful to compare this answer to the result from Example 9.4.1. Note that despite the similarity of the mean and standard deviation, the 95% confidence interval based on the t-distribution is a bit wider than the confidence interval based on the normal distribution. The reason is that in Example 9.4.1 the standard deviation was known exactly, while in this example the standard deviation needed to be estimated from the sample. This introduces an additional source of random error into the problem and thus the confidence interval must be wider to ensure the same likelihood of containing the true value of μ . ■

EXAMPLE 9.4.3. Consider the example given at the start of the chapter. A coin is flipped 100 times with a result of 60 heads and 40 tails. What would be a 90% confidence interval for the actual probability p that the coin shows heads on any given flip?

We represent a flip by $X \sim \text{Bernoulli}(p)$. From the i.i.d. sample X_1, X_2, \dots, X_{100} we have $\hat{p} = 60/100 = 0.6$. Despite the fact that σ is unknown, it would be inappropriate to use a t -distribution for the confidence interval in this case because X is far from a normal distribution. But we may still appeal

to the Central Limit Theorem and accept the sample as providing a reasonable estimate for the standard deviation. That is, if $X \sim \text{Bernoulli}(0.6)$, then $\sigma = SD[X] = \sqrt{0.24}$. Using this approximation $\sigma \approx \sqrt{0.24}$ we may proceed as before..... ■

EXERCISES

Ex. 9.4.1. are t intervals always larger than Normal

9.5 HYPOTHESIS TESTING

The idea of hypothesis testing is another approach to comparing an observed quantity from a sample (such as \bar{X}) to an expected result based on an assumption about the distribution X . There are many different types of hypothesis tests. What follows is far from an exhaustive list, but we explore some particular forms of hypothesis testing built around four familiar random distributions – the z-test, the t-test, the F-test, and the χ^2 -test.

For any hypothesis test a “null hypothesis” is a specific conjecture made about the nature of the distribution X . This is compared to an “alternate hypothesis” that specifies a particular manner in which the null may be an inaccurate assumption.

A computation is then performed based on the differences between the sample data and the result which would have been expected if the null hypothesis were true. This computation results in a quantity called a “P-value” which describes the probability that sample would be at least as far from expectation as was actually observed. The nature of this comparison between observation and expectation varies according to the type of test performed and the assumptions of the null hypothesis. A small P-value is an indication that the sample would be highly unusual (casting doubt on the null hypothesis), while a large P-value indicates that the sample is quite consistent with the assumptions of the null.

9.5.1 The z-test: Test for sample mean when σ is known

Suppose $X \sim \text{Normal}(\mu, \sigma^2)$ where μ is an unknown mean, but σ is a known standard deviation. Let X_1, X_2, \dots, X_n be an i.i.d. sample from the distribution X . Select as a null hypothesis the assumption that $\mu = c$ for some value c less than the observed average \bar{X} . Since the sample average \bar{X} is larger than the assumed mean, the assumption $\mu > c$ may be an appropriate alternate hypothesis. If the null is true, how likely is it we would have seen a sample mean as large as the observed value \bar{X} ?

To answer this question, we assume the empirical values of the sample X_1, X_2, \dots, X_n are known and let Y_1, Y_2, \dots, Y_n be an i.i.d. sample from the same distribution X . The Y_j variables effectively mimic the sampling procedure, an idea that will be commonly used through all tests of significance we consider.

We then calculate $P(\bar{Y} \geq \bar{X})$ where \bar{Y} is viewed as a random variable and \bar{X} is taken as the (deterministic) observed sample average. The \bar{X} statistic calculated from the observed data is known as the “test statistic”. The probability $P(\bar{Y} \geq \bar{X})$ describes how likely it is that the test statistic would be at least as far away from μ as what was observed. This probability can be computed precisely because the distribution of \bar{Y} is known exactly. Specifically, if the null hypothesis is true, then

$$Z = \frac{\sqrt{n}(\bar{Y} - c)}{\sigma} \sim \text{Normal}(0, 1).$$

So

$$\begin{aligned} P(\bar{Y} \geq \bar{X}) &= P\left(\frac{\sqrt{n}(\bar{Y} - \mu)}{\sigma} \geq \frac{\sqrt{n}(\bar{X} - c)}{\sigma}\right) \\ &= P(Z \geq \frac{\sqrt{n}(\bar{X} - c)}{\sigma}). \end{aligned}$$

In practice, before a sample is taken, a “significance level” $\alpha \in (0, 1)$ is typically selected. If $P(\bar{Y} \geq \bar{X}) < \alpha$ then the sample average is so far from the assumed mean c that the assumption $\mu = c$ is judged to

be incorrect. This is known as “rejecting the null hypothesis”. Alternatively if $P(\bar{Y} \geq \bar{X}) \geq \alpha$ then the sample mean is seen as consistent with the assumption $\mu = c$ and the null hypothesis is not rejected.

EXAMPLE 9.5.1. Suppose X has a normal distribution with known standard deviation $\sigma = 3.0$. A sample X_1, X_2, \dots, X_{16} of i.i.d. random variables is taken, each with distribution X . If the observed sample mean is $\bar{X} = 10.2$, what conclusion would a z-test reach if the null hypothesis assumes $\mu = 9.5$ (against an alternate hypothesis $\mu > 9.5$) at a significance level of $\alpha = 0.05$? ■

Under the null hypothesis, $Y_1, Y_2, Y_3, \dots, Y_{16} \sim \text{Normal}(9.5, 3.0)$ and independent, so

$$\begin{aligned} P(\bar{Y} \geq \bar{X}) &= P(\bar{Y} \geq 10.2) \\ &= P\left(\frac{4(\bar{Y} - 9.5)}{3} \geq \frac{4(10.2 - 9.5)}{3}\right) \\ &= P(Z \geq \frac{4(10.2 - 9.5)}{3}) \approx 0.175 \end{aligned}$$

where the final approximation is made using the fact that $Z \sim \text{Normal}(0, 1)$. The 0.175 figure is the P-value. Since it is larger than our significance level of $\alpha = 0.05$ we would not reject null hypothesis. Put another way, if the $\mu = 9.5$ assumption is true, the sampling procedure will produce a result at least as large as the sample average $\bar{X} = 10.2$ about 17.5% of the time. This is common enough that we cannot reject the $\mu = 9.5$ assumption.

EXAMPLE 9.5.2. Make the same assumptions as in Example 9.5.1, but this time test a null hypothesis that $\mu = 8.5$ (with an alternate hypothesis $\mu > 8.5$ and a significance level of $\alpha = 0.05$). Under the null hypothesis $Y_1, Y_2, Y_3, \dots, Y_{16} \sim \text{Normal}(8.5, 3.0)$ and are independent, so

$$\begin{aligned} P(\bar{Y} \geq \bar{X}) &= P(\bar{Y} \geq 10.2) \\ &= P(Z \geq \frac{4(10.2 - 8.5)}{3}) \approx 0.012 \end{aligned}$$

Since 0.012 is less than $\alpha = 0.05$ the null hypothesis would be rejected and the test would reach the conclusion that the true mean μ is some value larger than 8.5. Put another way, if the $\mu = 8.5$ assumption is true, the sampling procedure will produce a result as large as the sample average $\bar{X} = 10.2$ only about 1.2% of the time. This is rare enough that we can reject the hypothesis that $\mu = 8.5$. ■

For a large sample, a z-test is commonly used even without the assumption that X has a normally distribution. This is justified by appealing to the Central Limit Theorem.

EXAMPLE 9.5.3. Suppose a programmer is writing an app to identify faces based on digital photographs taken from social media. She wants to be sure that the app makes an accurate identification more than 90% of the time in the long run. She takes a random sample of 500 such photos and her app makes the correct identification 462 times - a 92.4% success rate. The programmer is hoping that this is an indication her app has a better than 90% success rate in the long run. However it is also possible the long term success rate is only 90%, but that the app happened to overperform this bar on the 500 photo sample. What does a z-test say about a null hypothesis that the app is only 90% accurate (compared to an alternate hypothesis that the app is more than 90% accurate with a significance level of $\alpha = 0.05$)?

The random variables in question are modeled by a Bernoulli distribution, as the app either reaches a correct conclusion or it does not. Under the null hypothesis $Y_1, Y_2, \dots, Y_{500} \sim \text{Bernoulli}(0.9)$ and are independent. The sample proportion does not precisely have a normal distribution, but the Central Limit Theorem implies that the standardized quantity

$$\frac{Y_1 + \dots + Y_{500} - 450}{\sqrt{45}}$$

should have approximately a $\text{Normal}(0, 1)$ distribution. Therefore

$$\begin{aligned} P\left(\frac{Y_1 + \dots + Y_{500}}{500} \geq \frac{X_1 + \dots + X_{500}}{500}\right) &\geq P\left(\frac{Y_1 + \dots + Y_{500} - 450}{\sqrt{45}} \geq \frac{X_1 + \dots + X_{500} - 450}{\sqrt{45}}\right) \\ &\approx P\left(Z \geq \frac{462 - 450}{\sqrt{45}}\right) \approx 0.03 \end{aligned}$$

Since 0.03 is less than $\alpha = 0.05$ the null hypothesis would be rejected and the test would reach the conclusion that the success rate for the app is greater than 90% . ■

The examples above concern tests on the right hand tail of a normal curve. That is, they test a null hypothesis $\mu = c$ against an alternate hypothesis $\mu > c$. It is also possible to perform a test on the left hand tail (testing a null hypothesis $\mu = c$ against an alternate hypothesis $\mu < c$) and even a two-tailed test (testing a null $\mu = c$ against an alternate $\mu \neq c$), an example of which follows below.

EXAMPLE 9.5.4. Suppose X has a normal distribution random variable with unknown mean and $\sigma = 6$. Suppose X_1, X_2, \dots, X_{25} is an i.i.d. sample taken with distribution X and that $\bar{X} = 6.2$. What conclusion would a z-test reach if the null hypothesis assumes $\mu = 4$ against an alternate hypothesis $\mu \neq 4$ at a significance level of $\alpha = 0.05$? Since the alternate hypothesis doesn't specify a particular direction in which the null may be incorrect, the appropriate probability to compute is

$$P(|\bar{Y} - 4| \geq |\bar{X} - 4|),$$

the probability that the absolute distance of a sample from the anticipated mean of 4.0 is larger than what was actually observed.

$$\begin{aligned} P(|\bar{Y} - 4| \geq |\bar{X} - 4|) &= 1 - P(|\bar{Y} - 4| < 2.2) \\ &= 1 - P(-2.2 \leq \bar{Y} - 4 \leq 2.2) \\ &= 1 - P\left(\frac{5(-2.2)}{6} \leq \frac{5(\bar{Y} - 4)}{6} \leq \frac{5(2.2)}{6}\right) \\ &= 1 - P\left(\frac{-11}{6} \leq Z \leq \frac{11}{6}\right) \approx 0.0668 \end{aligned}$$

since $Z \sim Normal(0, 1)$. As 0.0668 is slightly above the required significance level $\alpha = 0.05$ the test would not reject the null hypothesis. ■

9.5.2 The t-test: Test for sample mean when σ is unknown

As in the case of confidence intervals, when σ is unknown and estimated from the sample standard deviation S , an adjustment must be made by using the t-distribution.

Suppose X is known to be normally distributed with $X \sim Normal(\mu, \sigma^2)$ where μ and σ are unknown. Let X_1, X_2, \dots, X_n be an i.i.d. sample from the distribution X . Select as a null hypothesis that $\mu = c$ and select $\mu > c$ as an alternate hypothesis. Regard X_1, \dots, X_n as empirically known and let Y_1, \dots, Y_n be i.i.d. random variables which mimic the sampling procedure.

Under the null $Y_1, \dots, Y_n \sim Normal(c, \sigma^2)$ and are independent, from Corollary 8.1.10,

$$\frac{\sqrt{n}(\bar{Y} - c)}{S} \sim t_{n-1}$$

and so

$$\begin{aligned} P(\bar{Y} \geq \bar{X}) &= P\left(\frac{\sqrt{n}(\bar{Y} - c)}{S} \geq \frac{\sqrt{n}(\bar{X} - c)}{S}\right) \\ &= P(T \geq \frac{\sqrt{n}(\bar{X} - c)}{S}). \end{aligned}$$

The other aspects of the hypothesis test are the same except that the t-distribution must be used to calculate this final probability. As with the z-test, this could be performed as a one-tailed or a two-tailed test depending on the appropriate alternate hypothesis.

EXAMPLE 9.5.5. Suppose X has a normal distribution with unknown standard deviation. A sample X_1, X_2, \dots, X_{16} of i.i.d. random variables is taken, each with distribution X . The sample standard deviation is $S = 3.0$. What conclusion would a t-test reach if the null hypothesis assumes $\mu = 9.5$ at a significance level of $\alpha = 0.05$ (against the alternative hypothesis that $\mu > 9.5$)? ■

Under the null hypothesis, $Y_1, Y_2, \dots, Y_{16} \sim \text{Normal}(9.5, \sigma)$ and independent, so

$$\begin{aligned} P(\bar{Y} \geq \bar{X}) &= P(\bar{Y} \geq 10.2) \\ &= P\left(\frac{4(\bar{Y} - 9.5)}{3} \geq \frac{4(10.2 - 9.5)}{3}\right) \\ &= P\left(T \geq \frac{4(10.2 - 9.5)}{3}\right) \\ &= P\left(T \geq \frac{14}{15}\right) \approx 0.183 \end{aligned}$$

since $T \sim t_{15}$. As $0.183 > \alpha = 0.05$ the null hypothesis would not be rejected.

It is informative to compare this to Example 9.5.1 which was identical except that it was assumed that $\sigma = 3.0$ was known exactly rather than estimated from the sample. Note that the P-value in the case of the t-test (0.183) was slightly larger than in the case of the z-test (0.175). The reason is the use of S , a random variable, in place of σ , a deterministic constant. This adds an additional random factor into the computation and therefore makes larger deviations from the mean somewhat more likely.

9.5.3 A critical value approach

An alternate way to view the tests above is to focus on a “critical value”. Such a value is the dividing line beyond which the null will be rejected. If a test is being performed with a significance level α , then we can determine ahead of time where this line is and immediately reach a conclusion from the value of the test statistic without calculating a P-value. To demonstrate this we will redo Example 9.5.1 using this approach. In that example, X had a normal distribution with known standard deviation $\sigma = 3.0$ and Y_1, Y_2, \dots, Y_{16} were i.i.d. with distribution X . The null hypothesis assumed $\mu = 9.5$ while the alternate hypothesis was $\mu > 9.5$. The test had a significance level of $\alpha = 0.05$.

To find the critical value, we begin with the same computation as in Example 9.5.1, but keep \bar{X} as a variable.

$$\begin{aligned} P(\bar{Y} \geq \bar{X}) &= P\left(\frac{4(\bar{Y} - 9.5)}{3} \geq \frac{4(\bar{X} - 9.5)}{3}\right) \\ &= P(Z \geq \frac{4(\bar{X} - 9.5)}{3}) \end{aligned}$$

Whether or not the null is rejected depends entirely on whether this probability is above or below the significance level $\alpha = 0.05$, so the relevant question is what value of c ensures $P(Z \geq c) = 0.05$. This is something that can be calculated using Normal tables Table D.2 and in fact, $c \approx 1.645$. We solve the equation

$$\frac{4(\bar{X} - 9.5)}{3} = 1.645$$

which yields $\bar{X} \approx 10.73$. The figure 10.73 is the critical value. It is the dividing line we were seeking; if the sample average is above 10.73 the null hypothesis will be rejected while a sample average less than 10.73 will not cause the test to reject the null. For this particular example it was assumed that the observed sample average was $\bar{X} = 10.2$ which is why the null hypothesis was not rejected.

9.5.4 The χ^2 -test : Test for sample variance

Suppose instead of an inquiry about an average, we are interested in the variability in a population. Suppose $X \sim \text{Normal}(\mu, \sigma^2)$ with unknown σ . As with the previous hypothesis tests, we assume the empirical values of the sample X_1, X_2, \dots, X_n are known, and let Y_1, Y_2, \dots, Y_n be an i.i.d. sample with distribution X which mimics the sampling procedure. Select as a null hypothesis the assumption that $\sigma = c$ (and take $\sigma > c$ as an alternate hypothesis). How likely is it the sampling procedure would produce a sample standard deviation as large as S_X ?

We wish to calculate $P(S_Y \geq S_X)$, the probability a sample would produce a standard deviation at least as large as what was observed. Under the null hypothesis, from Theorem 8.1.9

$$\frac{(n-1)}{c^2} S_Y^2 \sim \chi_{n-1}^2.$$

Therefore,

$$\begin{aligned} P(S_Y \geq S_X) &= P(S_Y^2 \geq S_X^2) \\ &= P\left(\frac{(n-1)}{c^2} S_Y^2 \geq \frac{(n-1)}{c^2} S_X^2\right) \\ &= P(W \geq \frac{(n-1)}{c^2} S_X^2) \end{aligned}$$

which may be calculated since $W \sim \chi_{n-1}^2$ and n , c , and S_X are known.

EXAMPLE 9.5.6. Suppose X is normally distributed with unknown standard deviation σ . Let X_1, X_2, \dots, X_{16} be an i.i.d. sample with distribution X and a sample standard deviation $S_X = 3.5$. What conclusion would a χ^2 -test reach if the null hypothesis assumes $\sigma = 3$, with an alternate hypothesis that $\sigma > 3$, and a significance level of $\alpha = 0.05$?

The null hypothesis ensures

$$\begin{aligned} P(S_Y \geq S_X) &= P\left(\frac{15}{\sigma^2} S_Y^2 \geq \frac{15}{\sigma^2} S_X^2\right) \\ &= P(W \geq \frac{15}{9}(3.5)^2) \approx 0.157 \end{aligned}$$

since $W \sim \chi_{15}^2$. So there is about a 15.7% chance that a sample of size 16 would produce a standard deviation as large as 3.5. As this P-value is larger than α the null hypothesis is not rejected. ■

The shapes of the normal distribution and the t-distribution are symmetric about their means. This implies that when considering an interval centered at the mean, the two tail probabilities are always equal to each other. For example, if $Z \sim Normal(\mu, \sigma^2)$, then $P(Z \geq \mu + c) = P(Z \leq \mu - c)$ regardless of the value of c . In particular, when carrying out a computation for a hypothesis test,

$$\begin{aligned} P(|Z - \mu| \geq c) &= P(Z \geq \mu + c) + P(Z \leq \mu - c) \\ &= 2P(Z \geq \mu + c) \end{aligned}$$

since both tails have the same probability. However, this is not so for the chi-squared distribution. When performing a two-tailed test involving a distribution which is not symmetric, the interval selected is the one which has equal tail probabilities, each of which equal half of the confidence level. Due to this fact it is usually best to use a critical value approach.

EXAMPLE 9.5.7. Suppose X is normally distributed with unknown standard deviation σ . Let X_1, X_2, \dots, X_{16} be an i.i.d. sample with distribution X and a sample standard deviation $S_X = 3.5$, under the null hypothesis. What conclusion would a χ^2 -test reach if the null hypothesis assumes $\sigma = 3$, with an alternate hypothesis that $\sigma \neq 3$, and a significance level of $\alpha = 0.05$?

As in the previous example, we let Y_1, \dots, Y_{16} replicate the sampling procedure and use that $\frac{15}{\sigma^2} S_Y^2$ has a χ_{15}^2 distribution. With a significance level of $\alpha = 0.05$ the critical points will be the values of the χ_{15}^2 distribution that correspond to tail probabilities of 2.5%. It may be calculated that if $W \sim \chi_{15}^2$, then $P(W \leq 6.26) \approx 0.025$ while $P(W \geq 27.49) \approx 0.025$. As is readily observed 6.26 is the 0.025-th quantile and 27.49 is the 0.975-th quantile. They define the low and high critical values beyond which would be considered among the 5% most unusual occurrences for W . The corresponding observed value in the sample is $\frac{15}{9}(3.5)^2 \approx 20.42$ does not put it among this unusual 5% mark, so the null hypothesis would not be rejected. ■

9.5.5 The two-sample z-test: Test to compare sample means

Hypothesis tests may also be used to compare two samples to each other to see if the populations they were derived from were similar. This is of particular use in many applications. For instance: are the political opinions of one region different from another?; or are test scores at one school better than those at another school? These questions could be approached by taking random samples from each population and comparing them with each other.

Suppose X_1, X_2, \dots, X_{n_1} is an i.i.d. sample from a distribution $X \sim \text{Normal}(\mu_1, \sigma_1^2)$ and suppose Y_1, Y_2, \dots, Y_{n_2} is an i.i.d. sample from a distribution $Y \sim \text{Normal}(\mu_2, \sigma_2^2)$ independent of the X_j variables. Assume σ_1 and σ_2 are known, but μ_1 and μ_2 are not. How might we test a null hypothesis that $\mu_1 = \mu_2$ against an alternative hypothesis $\mu_1 \neq \mu_2$?

If the null hypothesis were true $\mu_1 - \mu_2 = 0$. We could calculate $\bar{X} - \bar{Y}$ and determine if this difference was close enough to 0 to make the null plausible. As usual, we mimic the sampling procedure, this time with both samples. Let V_1, \dots, V_{n_1} be an i.i.d. sample with distribution X and let W_1, \dots, W_{n_2} be an i.i.d. sample with distribution Y independent of the V_j variables. We would then calculate

$$P(|\bar{V} - \bar{W}| \geq |\bar{X} - \bar{Y}|),$$

the probability that the difference of sample averages would be at least as large as what was observed.

As $\bar{V} \sim \text{Normal}(\mu_1, \frac{\sigma_1^2}{n_1})$ and $\bar{W} \sim \text{Normal}(\mu_2, \frac{\sigma_2^2}{n_2})$. Under the null hypothesis the mean of $\bar{V} - \bar{W}$ is zero and they are independent with each having normal distribution. By Theorem 6.3.13 $\bar{V} - \bar{W} \sim \text{Normal}(0, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$. Therefore,

$$\begin{aligned} P(|\bar{V} - \bar{W}| \geq |\bar{X} - \bar{Y}|) &= P\left(\left|\frac{\bar{V} - \bar{W}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right| \geq \frac{|\bar{X} - \bar{Y}|}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right) \\ &= P(|Z| \geq \frac{|\bar{X} - \bar{Y}|}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}) \\ &= 2P(Z \leq -\frac{|\bar{X} - \bar{Y}|}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}) \end{aligned}$$

where $Z \sim \text{Normal}(0, 1)$.

EXAMPLE 9.5.8. Suppose a biologist wants to know if the average weights of adult members of a species of squirrel in one forest is the same as an identical species in a different location. Historical data suggests that the weights of the species have a standard deviation $\sigma = 10$ grams and the biologist is willing to use this assumption for his computations. Suppose he takes a sample of 30 squirrels from each location is is willing to regard these as independent i.i.d. samples from the respective populations. The first sample average is 122.4 grams and the second sample average is 127.6 grams. What conclusion would a two-sample z-test reach testing a null hypothesis that the population averages are the same against a alternate hypothesis that they are different at a significance level of $\alpha = 0.05$?

Nothing in the statement of the problem suggests that an assumption about the normality of the populations, so we will need to appeal to the Central Limit Theorem and be content that this is a decent approximation. Let X and Y represent the distributions of weights of the two populations. Under the null hypothesis these distributions have equal means $\mu = \mu_X = \mu_Y$ and $\sigma = 10$ for both distributions. Let V_1, \dots, V_{30} and W_1, \dots, W_{30} be i.i.d samples from populations X and Y respectively. Observe that

$$\bar{V} - \bar{W} = \frac{1}{30} \sum_{i=1}^{30} (V_i - W_i).$$

Now $V_1 - W_1, V_2 - W_2, V_3 - W_3, \dots, V_{30} - W_{30}$ are i.i.d. with zero mean and standard deviation $\sqrt{10^2 + 10^2} = 10\sqrt{2}$. By the Central Limit Theorem, the distribution of

$$\frac{\sqrt{30}(\bar{V} - \bar{W})}{10\sqrt{2}}$$

is approximately standard normal. Therefore,

$$\begin{aligned} P(|\bar{V} - \bar{W}| \geq 5.2) &= P\left(\left|\frac{\sqrt{30}(\bar{V} - \bar{W})}{10\sqrt{2}}\right| \geq \frac{\sqrt{30}(5.2)}{10\sqrt{2}}\right) \\ &= P(Z \geq 0.52\sqrt{15}) \approx 0.0440 \end{aligned}$$

where $Z \sim \text{Normal}(0, 1)$.

Since the P-value falls below the significance level, we would reject the null hypothesis and conclude that the populations have different average weights. ■

9.5.6 The F-test: Test to compare sample variances.

Let X_1, X_2, \dots, X_{n_1} be an i.i.d. sample from a distribution $X \sim \text{Normal}(\mu_1, \sigma_1^2)$ and Y_1, Y_2, \dots, Y_{n_2} be an i.i.d. sample (independent of the X_j variables) from a distribution $Y \sim \text{Normal}(\mu_2, \sigma_2^2)$. Suppose we wish to test the null hypothesis $\sigma_1 = \sigma_2$ against the alternate hypothesis $\sigma_1 \neq \sigma_2$.

Let V_1, \dots, V_{n_1} replicate the sample from X and W_1, \dots, W_{n_2} replicate the sample from Y . Let S_V^2 and S_W^2 denote the respective sample variances. Based on the previous examples it may be tempting to perform a test based on the probability that $|S_V^2 - S_W^2|$ is as large as the observed difference $|S_X^2 - S_Y^2|$. However, the random variable $S_V^2 - S_W^2$ does not have a distribution we have already considered. Instead, we will look at the ratio $\frac{S_V^2}{S_W^2}$. If the null hypothesis is true, this value should be close to 1 and the ratio has the benefit of being related to a familiar F -distribution.

The random variables $\frac{(n_1-1)}{\sigma_1^2}S_V^2$ and $\frac{(n_2-1)}{\sigma_2^2}S_W^2$ are independent and by Theorem 8.1.9 have the distributions $\chi_{n_1-1}^2$ and $\chi_{n_2-1}^2$ respectively. Therefore from Example 8.1.7 the ratio $\frac{S_V^2 \cdot \sigma_2^2}{S_W^2 \cdot \sigma_1^2}$ has a $F(n_1 - 1, n_2 - 1)$ distribution, and under the null hypothesis, this ratio simplifies to

$$\frac{S_V^2}{S_W^2} \sim F(n_1 - 1, n_2 - 1).$$

Since this is a distribution for which we may compute associated probabilities, we may use it to perform the hypothesis test. As the F distribution is not symmetric, we take a critical values approach.

EXAMPLE 9.5.9. Suppose X_1, X_2, \dots, X_{30} is an i.i.d. sample from a distribution $X \sim \text{Normal}(\mu_1, \sigma_1^2)$ and suppose Y_1, Y_2, \dots, Y_{25} is an i.i.d. sample from a distribution $Y \sim \text{Normal}(\mu_2, \sigma_2^2)$ independent of the X_j variables. If $S_X^2 = 11.4$ and $S_Y^2 = 5.1$, what conclusion would an F -test reach for null hypothesis suggesting $\sigma_1 = \sigma_2$, an alternate hypothesis suggesting $\sigma_1 \neq \sigma_2$, and a significance level of $\alpha = 0.05$?

From the computation above, if $R = \frac{S_V^2}{S_W^2}$ is the ratio of the sample variances, then $R \sim F(29, 24)$. From this, it may be calculated that if $P(R \leq 0.464) \approx 0.025$ while $P(R \geq 2.22) \approx 0.025$. Since the observed ratio $\frac{11.4}{5.1} \approx 2.24$ is outside of the interval $(0.464, 2.22)$, the null hypothesis would be rejected. ■

9.5.7 A χ^2 -test for “goodness of fit”

A more common use of the χ^2 is for something called a “goodness of fit” test. In this case we seek to determine whether the distribution of results in a sample could plausibly have come from a distribution specified by a null hypothesis. The test statistic is calculated by comparing the observed count of data points within specified categories relative to the expected number of results in those categories according to the assumed null.

Specifically, let X is a random variable with finite range $\{c_1, c_2, \dots, c_k\}$ for which $P(X = c_j) = p_j > 0$ for $1 \leq j \leq k$. Let X_1, X_2, \dots, X_n be the empirical results of a sample from the distribution X and let

$Y_j = |\{j : X_j = c_j\}|$. That is, Y_j is the number of data points in the sample that resulted in the c_j outcome. Then the statistic $\chi^2 = \sum_{j=1}^k \frac{(Y_j - np_j)^2}{np_j}$ has approximately the $\chi^2(k-1)$ distribution. Note that the np_j term is the expected number of observations of the outcome c_j , so the numerator of the fractions in the χ^2 computation measures the squared difference between the observed and expected counts.

A formal proof of the statement above requires a deeper understanding of linear algebra than we assume as a prerequisite for this text, but below we demonstrate its truth in the special case where $k = 2$ and we provide the formal technical details in the appendix. The approximation itself relies on the Central Limit Theorem and as with that theorem, larger values of n will tend to lead to a better approximation. Before proceeding to the proof, we present an example will help illustrate the use of the test.

EXAMPLE 9.5.10. The student body at an undergraduate university is 20% seniors, 24% juniors, 26% sophomores, and 30% freshman. Suppose a researcher takes a sample of 50 such students. Within the sample there are 13 seniors, 16 juniors, 10 sophomores, and 11 freshmen. The researcher claims that his sampling procedure should have produced independent selections from the student body, with each student equally likely to be selected. Is this a plausible claim given the observed results?

If the claim is true (which we take as the null hypothesis), then selecting an individual for the sample should be like the empirical result from a random variable X with a range of { senior, junior, sophomore, freshman } where the probabilities of each outcome are the percentages described above. For instance, $p_{senior} = P(X = "senior") = 0.2$. The expected value of results based on the null hypothesis is then 10 seniors, 12 juniors, 13 sophomores, and 15 freshmen. So,

$$\chi^2 = \frac{(13-10)^2}{10} + \frac{(16-12)^2}{12} + \frac{(10-13)^2}{13} + \frac{(11-15)^2}{15} \approx 3.99$$

Notice that the χ^2 statistic can never be less than zero, and that when observed results are close to what was expected, the resulting fraction is small and does not contribute much to the sum. It is only when there is a relatively large discrepancy between observation and expectation that χ^2 will have a large value.

Since there were four categories, if the null hypothesis is correct, this statistic resulted from a $\chi^2(3)$ distribution. To see if such a thing is plausible, we let $W \sim \chi^2(3)$ and calculate that $P(W \geq 3.99) \approx 0.2625$. The researcher's claim seems plausible. According to the χ^2 -test, samples this far from expectation should be observed around 26% of the time. ■

To give a bit more insight into this test, consider each term individually. Each of the variables Y_1, Y_2, \dots, Y_n is a binomial random variable. For example Y_1 represents the number of times the outcome c_1 is observed in n trials, when $P(X = c_1) = p_1$, so $Y_1 \sim \text{Binomial}(n, p_1)$. Therefore $E[Y_1] = np_1$ and $SD[Y_1] = \sqrt{np_1(1-p_1)}$. From the Central Limiit Theorem, the normalized quantity $\frac{Y_1 - np_1}{\sqrt{np_1(1-p_1)}}$ has approximately the Normal(0, 1) distribution, and therefore its square $\frac{(Y_1 - np_1)^2}{np_1(1-p_1)}$ has approximately a $\chi^2(1)$ distribution. Except for the $(1 - p_1)$ term in the denominator, this is the first fraction in the sum of our test statistic. The additional factor in the denominator is connected to the reason the resulting distribution has $n - 1$ degrees of freedom instead of n degrees of freedom; the variables Y_1, Y_2, \dots, Y_n are dependent. Untangling this dependence in the general case is complicated, but if $k = 2$ we can prove a rigorous statement without the use of linear algebra.

THEOREM 9.5.11. Let X be a random variable with finite range $\{c_1, c_2\}$ for which $P(X = c_j) = p_j > 0$ for $j = 1, 2$. Let X_1, X_2, \dots, X_n be an i.i.d. sample with distribution X and let $Y_j = |\{j : X_j = c_j\}|$. Then $\chi^2 = \sum_{j=1}^2 \frac{(Y_j - np_j)^2}{np_j}$ has the same distribution as $(\frac{Z - E[Z]}{SD[Z]})^2$ where $Z \sim \text{Binomial}(n, p_1)$

Proof - To simplify notation, let $p = p_1$ and note that $p_2 = 1 - p$. Also, let $N = Y_1$ and note $Y_2 = n - N$. Then,

$$\begin{aligned}\chi^2 &= \frac{(N - np)^2}{np} + \frac{((n - N) - n(1 - p))^2}{n(1 - p)} \\ &= \frac{N^2 - 2npN + n^2p^2}{np(1 - p)} \\ &= \left(\frac{N - np}{\sqrt{np(1 - p)}}\right)^2\end{aligned}$$

Since $N = Y_1 \sim \text{Binomial}(n, p)$, the result follows.

The central limit theorem guarantees $\frac{N - np}{\sqrt{np(1 - p)}}$ converges in distribution to a $\text{Normal}(0, 1)$ distribution, the χ^2 quantity will have approximately a $\chi^2(1)$ distribution for large values of n .

LINEAR REGRESSION

In Chapter 6 we discussed the concepts of covariance and correlation – two ways of measuring the extent to which two random variables, X and Y were related to each other. In many cases we would like to take this a step further and try to use information from one variable to make predictions about the outcome of the other. For instance

10.1 SAMPLE COVARIANCE AND CORRELATION

We have so far considered summarizing a set of observations where one measurement is made on each individual or unit, but often in real-life random experiments we make multiple measurements on each individual. For example, during a health check-up a doctor might record the height, weight, age, sex, pulse rate, and blood pressure.

Just as we did for single measurements, we can represent the observed data by their empirical distribution, which is now a function of multiple arguments. For example, if we measure two random variables (X_i, Y_i) for the i th individual (say weight and blood pressure), then the empirical distribution function is given by

$$f(t, s) = \frac{1}{n} \# \{X = t, Y = s\}.$$

We can now use this to estimate population features by the corresponding feature of the empirical distribution. For example, the population covariance $Cov[X, Y] = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$ gives a measure of how X and Y relate to each other. The sample version of this is the sample covariance

$$S_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \frac{1}{n} \sum_{i=1}^n X_i Y_i - \bar{X} \bar{Y}. \quad (10.1.1)$$

The sample correlation coefficient is defined similarly to population correlation coefficient $\rho[X, Y]$ as

$$r[X, Y] = \frac{S_{XY}}{S_X S_Y}, \quad (10.1.2)$$

where S_X and S_Y are the sample standard deviations of X and Y respectively. As with $\rho[X, Y]$, $r[X, Y]$ is bounded between -1 and 1 , and is invariant to scale and location transformations, that is, for real numbers a, b, c, d ,

$$r[aX + b, cY + d] = r[X, Y]$$

10.2 SIMPLE LINEAR MODEL

We will assume that the variable Y depends on X in a linear fashion, but that it is also affected by random factors. Specifically we will assume there is a regression line $y = \alpha + \beta x$ and that for given x-values X_1, X_2, \dots, X_n the corresponding y -values Y_1, Y_2, \dots, Y_n are given by

$$Y_j = \alpha + \beta X_j + \epsilon_j, \quad (10.2.1)$$

for $j = 1, 2, \dots, n$ and where each of the ϵ_j are independent random variables with $\epsilon_j \sim \text{Normal}(0, \sigma^2)$. Equation (10.2.1) is referred to as the simple linear model. In particular ϵ_j are the (random) vertical distance of the point (X_j, Y_j) from the regression line. For all results below we assume $\sigma^2 > 0$ is the variance of the errors, assumed to be the same for every data point. We also assume that not all of the X_j quantities are the same so that the variance of these quantities is non-zero. In particular this means $n \geq 2$.

10.3 THE LEAST SQUARES LINE

The values of $(X_1, Y_1), \dots, (X_n, Y_n)$ are collected data. Though we assume that this data is produced via the simple linear model, we typically do not know the actual values of the slope β or the y-intercept α . The goal of this section is to illustrate a way to estimate these values from the data.

For a line $y = a + bx$ the “residual” of a data point (X_j, Y_j) is defined to be the quantity $Y_j - (a + bX_j)$. This is the difference between the actual y-value of the data point and the location where the line predicts the y-value should be. In other words, it may be viewed as the error of the line when attempting to predict the y-value corresponding to the X_j data point. Among all possible lines through the data, there is one which minimizes the sum of these squared residual errors. This is called the “least squares line”.

Let $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ be points on the plane. Suppose we wish to find a line that minimizes the sum of squared residual errors. That is, let $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ be defined as

$$g(a, b) = \sum_{j=1}^n [Y_j - (a + bX_j)]^2.$$

The objective is to minimize g . So using calculus,

$$0 = \frac{\partial g}{\partial a} = -2 \sum_{j=1}^n [Y_j - a - bX_j] \quad (10.3.1)$$

and

$$0 = \frac{\partial g}{\partial b} = -2 \sum_{j=1}^n X_j [Y_j - a - bX_j]. \quad (10.3.2)$$

From equation (10.3.1) we have

$$0 = \sum_{j=1}^n [Y_j - a - bX_j] = \sum_{j=1}^n Y_j - \sum_{j=1}^n a - b \sum_{j=1}^n X_j = n\bar{Y} - na - bn\bar{X} = n(\bar{Y} - (a + b\bar{X}))$$

Therefore¹

$$\bar{Y} = a + b\bar{X}, \quad (10.3.3)$$

which shows that the point (\bar{X}, \bar{Y}) must lie on the least squares line. The point (\bar{X}, \bar{Y}) is known as the point of averages. Similarly from equation (10.3.2),

$$0 = \sum_{j=1}^n X_j [Y_j - a - bX_j] = \sum_{j=1}^n (X_j Y_j - aX_j - bX_j^2) = \sum_{j=1}^n X_j Y_j - an\bar{X} + b \sum_{j=1}^n X_j^2$$

so that

$$\sum_{j=1}^n X_j Y_j = an\bar{X} + b \sum_{j=1}^n X_j^2. \quad (10.3.4)$$

We now use the system of two equations (given by (10.3.3) and (10.3.4)) solve for a, b to get

$$b = \frac{\left(\sum_{j=1}^n X_j Y_j \right) - n\bar{X}\bar{Y}}{\left(\sum_{j=1}^n X_j^2 \right) - n\bar{X}^2} \quad (10.3.5)$$

$$(10.3.6)$$

¹ We shall use the notation $\bar{X}, \bar{Y}, S_X, S_Y, r[X, Y]$ (below), even though they are not necessarily random quantities. This is to simplify notation and will allow us to use known properties, in the event they are random.

Recall that the sample variance of X_1, X_2, \dots, X_n is

$$\begin{aligned} S_X^2 &= \frac{1}{n-1} \left[\sum_{j=1}^n (X_j - \bar{X})^2 \right] \\ &= \frac{1}{n-1} \left[\sum_{j=1}^n X_j^2 - 2\bar{X} \sum_{j=1}^n X_j + n\bar{X}^2 \right] \\ &= \frac{1}{n-1} \left[\left(\sum_{j=1}^n X_j^2 \right) - 2\bar{X} \left(\sum_{j=1}^n X_j \right) + n\bar{X}^2 \right] \\ &= \frac{1}{n-1} \left[\left(\sum_{j=1}^n X_j^2 \right) - 2n\bar{X}^2 + n\bar{X}^2 \right] \\ &= \frac{1}{n-1} \left[\left(\sum_{j=1}^n X_j^2 \right) - n\bar{X}^2 \right] \end{aligned}$$

Therefore, the denominator of (10.3.5) is simply $(n-1)S_X^2$. The numerator may be written more simply by using the notation of sample covariance and correlation defined in (10.1.1) and (10.1.2). So from (10.3.5) we have

$$b = \frac{\left(\sum_{j=1}^n X_j Y_j \right) - n\bar{X}\bar{Y}}{\left(\sum_{j=1}^n X_j^2 \right) - n\bar{X}^2} = \frac{(n-1)S_{XY}}{(n-1)S_X^2} = \frac{r[X, Y]S_Y}{S_X}$$

Using the above and (10.3.3), we also now can write a nice formula for a , which is

$$a = \bar{Y} - \frac{r[X, Y]S_Y}{S_X} \bar{X} \quad (10.3.7)$$

By the above calculation we have shown that the least squares line minimizing the sum of the squared residual errors is the line passing through the point of averages (\bar{X}, \bar{Y}) and having a slope equal to $b = \frac{r[X, Y]S_Y}{S_X}$. We state this precisely in the Theorem below.

THEOREM 10.3.1. *Let $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ be given data points. Then the least squares line passes through (\bar{X}, \bar{Y}) and has slope given by $\frac{r[X, Y]S_Y}{S_X}$.*

We illustrate the use of these formulas with two examples given below.

EXAMPLE 10.3.2. Consider the following five data points:

X	Y
3	6
4	5
5	6
6	4
7	2

These points are not colinear, but suppose we wish to find a line that most closely approximates their trend in the least squares sense described above. Viewing these as samples, it is routine to calculate that the formulas above yield $a = 9.1$ and $b = -0.9$. Of all of the lines in the plane, the one that minimizes the sum of squared residual errors for the data set above is the line $y = 9.1 - 0.9x$.

The R software also has a feature to perform a regression directly. To obtain this result using R we could first create vectors that represent the data:

```
> x <- c(3,4,5,6,7)
> y <- c(6,5,6,4,2)
```

And then instruct R to perform the regression using the command “lm” indicating the linear model.

```
> lm(y ~ x)
```

The order of the variables in this command is important with this $y \sim x$ indicating that the y variable is being predicted using the x variable as input.

The resulting output from R is

(Intercept)	x
9.1	-0.9

the values of the intercept and slope of the least squares line respectively. ■

EXAMPLE 10.3.3. Suppose as part of a health study, a researcher collects data for weights and heights of sixty adult men in a population. The average height of the men is 174 cm with a sample standard deviation of 8.0 cm. The average weight of the men is 78 kg with a sample standard deviation of 10 kg. The correlation between the variables in the sample was 0.55.

This information alone is enough to find the least squares line for predicting weight from height. The reader may use the formulas above to verify that $b = 0.6875$ and $a = -41.625$. Therefore, among all lines, $y = -41.625 + 0.6875x$ is the one which minimizes the sum of squared residuals.

This does not necessarily mean this line would be appropriate for predicting new data points. To make such a declaration, we would want to have some evidence that the two variables had a linear relationship to begin with, but regardless of whether or not the data was produced from a simple linear model, the line above minimizes error in the least squares sense. ■

EXERCISES

Ex. 10.3.1. Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be data produced via the simple linear model and suppose $y = a + bx$ is the least squares line for the data. Recall from above that the residual for any given data point is $Y_j - (a + bX_j)$, the error the line makes in predicting the correct y -value from the given x -value. Show that the sum of the residuals over all n data points must be zero.

Ex. 10.3.2. Suppose that instead of using the simple linear model, we assume the regression line is known to pass through the origin. That is, the regression line has the form $y = \beta x$ and for given x -values X_1, X_2, \dots, X_n the corresponding y -values Y_1, Y_2, \dots, Y_n are given by

$$Y_j = \beta X_j + \epsilon_j, \tag{10.3.8}$$

for $j = 1, 2, \dots, n$. As with the simple linear model, we assume each of the ϵ_j are independent random variables with $\epsilon_j \sim \text{Normal}(0, \sigma^2)$. (We will refer to this as the “linear model though the origin” and will have several exercises investigating how several formulas from this chapter would need to be modified for such a model.)

Assuming data $(X_1, Y_1), \dots, (X_n, Y_n)$ was produced from the linear model through the origin, find the least squares line through the origin. That is, find a formula for b such that the line $y = bx$ minimizes the sum of squared residual errors.

10.4 a AND b AS RANDOM VARIABLES

In this section (and the remainder of this chapter) we will assume that $(X_1, Y_1), \dots, (X_n, Y_n)$ follow the simple linear model (10.2.1). In other words, there is a regression line $y = \alpha + \beta x$ and that for given x -values X_1, X_2, \dots, X_n the corresponding y -values Y_1, Y_2, \dots, Y_n are given by (10.2.1). In the previous section this data was used to produce a mean squared error-minimizing least squares line $y = a + bx$. In this section we investigate how well the random quantities a and b approximate the (unknown) values α and β .

THEOREM 10.4.1. Under the assumptions of the simple linear model (10.2.1), the slope b of the least squares line is a linear combination of the Y_j variables. Further it has a normal distribution with mean β and variance $\frac{\sigma^2}{(n-1)S_X^2}$.

Proof - First recall that the $X_1, X_2, X_3, \dots, X_n$ are assumed to be deterministic, so will be treated as known constants. The data points Y_1, Y_2, \dots, Y_n are assumed to follow the simple linear model (10.2.1). So for $j = 1, \dots, n$,

$$\begin{aligned} E[Y_j] &= E[\alpha + \beta X_j + \epsilon_j] = \alpha + \beta X_j + E[\epsilon_j] = \alpha + \beta X_j \\ &\text{and} \\ Var[Y_j] &= Var[\alpha + \beta X_j + \epsilon_j] = Var[\epsilon_j] = \sigma^2. \end{aligned}$$

Using the formula, (10.3.5), we derived for b and the above we have

$$\begin{aligned} E[b] &= E \left[\frac{(\sum_{j=1}^n X_j Y_j) - n\bar{X}\bar{Y}}{(n-1)S_X^2} \right] \\ &= \frac{1}{(n-1)S_X^2} \left[(\sum_{j=1}^n X_j E[Y_j]) - n\bar{X}E[\bar{Y}] \right] \\ &= \frac{1}{(n-1)S_X^2} \left[(\sum_{j=1}^n X_j (\alpha + \beta X_j)) - n\bar{X}(\alpha + \beta \bar{X}) \right] \\ &= \frac{1}{(n-1)S_X^2} \left[n\alpha\bar{X} + \beta(\sum_{j=1}^n X_j^2) - n\alpha\bar{X} - \beta n\bar{X}^2 \right] \\ &= \frac{\beta}{(n-1)S_X^2} \left[(\sum_{j=1}^n X_j^2) - n\bar{X}^2 \right] = \beta. \end{aligned}$$

Similarly,

$$\begin{aligned} Var[b] &= Var \left[\frac{(\sum_{j=1}^n X_j Y_j) - n\bar{X}\bar{Y}}{(n-1)S_X^2} \right] \\ &= \frac{1}{[(n-1)S_X^2]^2} \left[(\sum_{j=1}^n X_j^2 Var[Y_j]) - n^2\bar{X}^2 Var[\bar{Y}] \right] \\ &= \frac{1}{[(n-1)S_X^2]^2} \left[(\sum_{j=1}^n X_j^2 \sigma^2) - n^2\bar{X}^2 (\sigma^2/n) \right] \\ &= \frac{\sigma^2}{[(n-1)S_X^2]^2} \left[(\sum_{j=1}^n X_j^2) - n\bar{X}^2 \right] \\ &= \frac{\sigma^2}{(n-1)S_X^2}. \end{aligned}$$

The algebra below justifies that b is a linear combination of the Y_j variables.

$$b = \frac{(\sum_{j=1}^n X_j Y_j) - n\bar{X}\bar{Y}}{(n-1)S_X^2} = \frac{1}{(n-1)S_X^2} \left[(\sum_{j=1}^n X_j Y_j) - (\sum_{j=1}^n \bar{X}Y_j) \right] = \sum_{j=1}^n \left[\frac{X_j - \bar{X}}{(n-1)S_X^2} \right] Y_j$$

Since b is a linear combination of independent, normal random variables Y_j , b itself is also a normal random variable (Theorem 6.3.13). ■

As noted above, the least squares line can be defined as the line of slope b passing through the point of averages. The following lemma is a useful fact about how these quantities relate to each other.

LEMMA 10.4.2. *Let b be the slope of the least squares line and let \bar{Y} be the sample average of the Y_j variables. Then b and \bar{Y} are independent.*

Proof - By Theorem 6.3.13, \bar{Y} has a normal distribution and so does b by Theorem 10.4.1. By Theorem 6.4.3, all we have to show is that \bar{Y} and b are uncorrelated. Note that the Y_j variables are all independent of each other and so $Cov[Y_j, Y_k]$ will be zero if $j \neq k$ and will equal the variance σ^2 otherwise. So,

$$\begin{aligned} Cov[b, \bar{Y}] &= Cov \left[\sum_{j=1}^n \frac{X_j - \bar{X}}{(n-1)S_X^2} Y_j, \frac{1}{n} \sum_{k=1}^n Y_k \right] \\ &= \sum_{j=1}^n \sum_{k=1}^n Cov \left[\frac{X_j - \bar{X}}{(n-1)S_X^2} Y_j, \frac{1}{n} Y_k \right] \\ &= \sum_{j=1}^n \sum_{k=1}^n \frac{X_j - \bar{X}}{n(n-1)S_X^2} Cov[Y_j, Y_k] \\ &= \sum_{j=1}^n \frac{X_j - \bar{X}}{n(n-1)S_X^2} \sigma^2 \\ &= \frac{\sigma^2}{n(n-1)S_X^2} \sum_{j=1}^n X_j - \bar{X} = 0. \end{aligned}$$
■

We conclude this section with a result on the distribution of a .

THEOREM 10.4.3. *Under the assumptions of the simple linear model (10.2.1), The y -intercept a (given by (10.3.7)) of the least squares line is a linear combination of Y_j variables. Further it has a normal distribution with mean α and variance $\sigma^2(\frac{1}{n} + \frac{\bar{X}^2}{(n-1)S_X^2})$.*

Proof- See Exercise 10.4.1.

EXERCISES

Ex. 10.4.1. Prove Theorem 10.4.3. (Hint: Make use of the fact that $\bar{Y} = a + b\bar{X}$ and what has previously been proven about \bar{Y} and b).

Ex. 10.4.2. Show that, generally speaking, a and b are not independent. Find necessary and sufficient conditions for when the two variables are independent.

Ex. 10.4.3. Show that a and \bar{Y} are never independent.

Ex. 10.4.4. Continuing from Exercise 10.3.2, assuming the regression line $y = \beta x$ passes through the origin and b is the least squares line of the form $y = bx$, do the following:

- (a) Find the expected value of b .

- (b) Find the variance of b .
- (c) Determine whether or not b has a normal distribution.
- (d) Determine if b and \bar{Y} are independent.

10.5 PREDICTING NEW DATA WHEN σ^2 IS KNOWN

In this section we return to question of using data for prediction. We continue to assume the simple linear model (10.2.1). We further assume that α and β are estimated by a and b (as calculated from the data $(X_1, Y_1), \dots, (X_n, Y_n)$) and parameter σ^2 describing the variability of data around the regression line is a known quantity.

First suppose for a particular deterministic x-value X^* that we want to use the data to estimate the corresponding y -value $Y^* = \alpha + \beta X^*$ on the regression line by $Y = a + bX^*$.

THEOREM 10.5.1. *The quantity $Y = a + bX^*$ has a normal distribution with mean $Y^* = \alpha + \beta X^*$ and variance $\sigma^2\left(\frac{1}{n} + \frac{(X^* - \bar{X})^2}{(n-1)S_X^2}\right)$.*

Proof - Recall from Theorem 10.4.3 and Theorem 10.4.1 that a and b are both linear combination of the random variables Y_j normal distribution. So Y has normal distribution by Theorem 6.3.13. We need to calculate only its mean and variance. The expected value is simple to calculate.

$$\begin{aligned} E[Y] &= E[a + bX^*] \\ &= E[a] + E[b]X^* \\ &= \alpha + \beta X^* = Y^* \end{aligned}$$

If a and b were independent, then calculating the variance of Y would also be a simple task, but this is typically this is not the case. However, from Lemma 10.4.2, we know that b and \bar{Y} are independent. To make use of this, using (10.3.3), we may rewrite the line in point-slope form around the point of averages: $Y = \bar{Y} + b(X^* - \bar{X})$. From this we have,

$$\begin{aligned} Var[Y] &= Var[\bar{Y} + b(X^* - \bar{X})] \\ &= Var[\bar{Y}] + Var[b](X^* - \bar{X})^2 \\ &= \frac{\sigma^2}{n} + \frac{\sigma^2}{(n-1)S_X^2}(X^* - \bar{X})^2 \\ &= \sigma^2 \left(\frac{1}{n} + \frac{(X^* - \bar{X})^2}{(n-1)S_X^2} \right). \end{aligned}$$

■

Note that for various values of X^* this variance is minimal when X^* is \bar{X} , the average value of the x-data. In this case $Var[Y] = \frac{\sigma^2}{n} = Var[\bar{Y}]$ as expected. The further X^* is from the average of the x-values, the more variance there is in predicting the point on the regression line.

Next suppose that, instead of trying to estimate a point on the regression line, we are trying to predict a new data point produced from the linear model. Let X^* now represent the x-value of some new data point and let $Y^* = \alpha + \beta X^* + \epsilon^*$ where $\epsilon^* \sim \text{Normal}(0, \sigma^2)$ where the random variable ϵ^* is assumed to be independent of all prior ϵ_j which produced the original data set. The following theorem addresses the distribution of the predictive error made when estimating Y^* by the quantity $Y = a + bX^*$.

THEOREM 10.5.2. *If (X^*, Y^*) is a new data point, as described in the previous paragraph, then the predictive error in estimating Y^* using the least square line is $(a + bX^*) - Y^*$ which is normally distributed with mean 0 and variance $\sigma^2\left(1 + \frac{1}{n} + \frac{(X^* - \bar{X})^2}{(n-1)S_X^2}\right)$.*

Proof - The expected value of the predictive error is zero since

$$\begin{aligned} E[(a + bX^*) - Y^*] &= E[a] + E[b]X^* - E[\alpha + \beta X^* + \epsilon^*] \\ &= \alpha + \beta X^* - \alpha - \beta X^* - E[\epsilon^*] = 0. \end{aligned}$$

Both quantities a and b are linear combinations of the Y_j variables and so

$$\begin{aligned} a + bX^* - Y^* &= a + bX^* - \alpha - \beta X^* - \epsilon^* \\ &= (-\alpha - \beta X^*) + \text{a linear combination of } Y_1, Y_2, \dots, Y_n, \epsilon^*. \end{aligned}$$

All $(n+1)$ of the variables, $Y_1, Y_2, \dots, Y_n, \epsilon^*$, are independent and have a normal distribution. As $(-\alpha - \beta X^*)$ is a constant, from the above $(a + bX^* - Y^*)$ has a normal distribution.

Finally, to calculate the variance, we again rewrite $a + bX^*$ in point-slope form and exploit independence.

$$\begin{aligned} \text{Var}[(a + bX^* - Y^*)] &= \text{Var}[(\bar{Y} + b(X^* - \bar{X}) - (\alpha + \beta X^* + \epsilon^*))] \\ &= \text{Var}[\bar{Y}] + \text{Var}[b](X^* - \bar{X})^2 + \text{Var}[\epsilon^*] \\ &= \frac{\sigma^2}{n} + \frac{\sigma^2}{(n-1)S_X^2}(X^* - \bar{X})^2 + \sigma^2 \\ &= \sigma^2 \left(1 + \frac{1}{n} + \frac{(X^* - \bar{X})^2}{(n-1)S_X^2} \right). \end{aligned}$$

EXAMPLE 10.5.3. A mathematics professor at a large university is studying the relationship between scores on a preparation assessment quiz students take on the first day of class and their actual percentage score at the end of class. Assuming the simple linear model with $\sigma = 6$, he takes a random sample of 30 students and discovers their average score on the quiz is $\bar{X} = 54$ with a sample standard deviation of $S_X = 12$, while the average percentage score in the class is $\bar{Y} = 68$ with a sample standard deviation of $S_Y = 10$. The sample correlation is $r[X, Y] = 0.6$. So according to the results above, the least squares line for predicting the course percentage from the preliminary quiz will be $y = 0.5x + 41$.

If we wish to use the line to predict the course percentage for someone who scores a 54 on the preliminary quiz, we would find $y = 0.5(54) + 41 = 68$, as expected since someone who gets an average score on the quiz is likely to get around the average percentage in the class.

Similarly if we wish to use the line to predict the course percentage for someone who scores a 80 on the preliminary quiz, we would find $y = 0.5(80) + 41 = 81$. Also not surprising. Due to the positive correlation, a student scoring above average on the quiz is also likely to score higher in the course as well.

The previous theorem allows us to go further and calculate a standard deviation associated with these estimates. For the student who scores a 54 on the preliminary quiz, let Y^* be the actual course percentage and let $a + bX^* = 68$ be the least squares line estimate we made above. Then,

$$\text{Var}[a + bX^* - Y^*] = 36\left(1 + \frac{1}{30} + 0\right) = 37.2$$

and so the standard deviation in the predictive error is $SD[a + bX^* - Y^*] \approx 6.1$. This means that students who make an average score of 54 on the preliminary quiz will have a range of percentages in the course. This range will have a normal distribution with mean 68 and standard deviation 6.1. We could then use normal curve computations to make further predictions about how likely such a student may be to reach a certain benchmark.

Next take the example of a student who scores 80 on the preliminary quiz. The least squares line predicts the course percentage for such a student will be $a + bX^* = 81$, but now

$$\text{Var}[a + bX^* - Y^*] = 36\left(1 + \frac{1}{30} + \frac{(80-54)^2}{29 \cdot 12^2}\right) \approx 43.0$$

and so $SD[a + bX^* - Y^*] \approx 6.6$. Student who score an 80 on the preliminary exam will have a range of course percentages with a normal distribution of mean 81 and standard deviation 6.6.

Thinking of the standard deviation as the likely error associated with prediction this example suggests that predictions of data further from the mean will tend to have less accuracy than predictions near to the mean. This is true in the simple linear model and will be explored in the exercises. ■

EXERCISES

Ex. 10.5.1. Using the figures from Example 10.5.3 do the following. Two students are selected independently at random. The first scored a 50 on the preliminary quiz while the second scored 60. Determine how likely it is that the student who scored the lower grade on the quiz will score a higher percentage in the course.

Ex. 10.5.2. Explain why $\text{Var}[a + bX^* - Y^*]$ is minimized when $X^* = \bar{X}$.

10.6 HYPOTHESIS TESTING AND REGRESSION

As a and b both have a normal distribution under the assumption of the simple linear model, it is possible to perform tests of significance concerning the values of α and β . Of particular importance is a test with a null hypothesis that $\beta = 0$ and an alternate hypothesis $\beta \neq 0$. This is commonly called a “test of utility”. The reason for this name is that if $\beta = 0$, then the simple linear model produces output values $Y_j = \alpha + \epsilon_j$ which do not depend on the corresponding input X_j . Therefore knowing the value of X_j should not be at all helpful in predicting the corresponding Y_j result. However, if $\beta \neq 0$ then knowing X_j should be at least somewhat useful in predicting Y_j value.

EXAMPLE 10.6.1. Suppose $(X_1, Y_1), \dots, (X_{16}, Y_{16})$ follows the simple linear model with $\sigma = 5$ and produces a least squares line $y = 0.3 + 1.1x$. Suppose the sample average of the X_j data is 20 and the sample variance is $S_X^2 = 10$. What is the conclusion of a test of utility at a significance level of $\alpha = 0.05$? ■

From the given least squares line, $b = 1.1$. As noted above, a test of utility compares a null hypothesis that $\beta = 0$ to an alternate hypothesis $\beta \neq 0$, so this will be a two-tailed test. If the null were true, then $E[b] = 0$ and we can use the normal distribution to determine whether the 1.1 value is so far from zero that the null seems unreasonable. Using the same sample mimicing idea introduced in Chapter 9 we let Z_1, \dots, Z_{16} be random variables produced from X_1, \dots, X_{16} via the simple linear model. From Theorem 10.4.1, the slope of the least squares line for the $(X_1, Z_1), \dots, (X_{16}, Z_{16})$ data has a normal distribution with mean $\beta = 0$ and variance $\frac{\sigma^2}{(n-1)S_X^2} = \frac{1}{6}$. Therefore we can calculate

$$\begin{aligned} P(|\text{slope of the least squares line}| \geq 1.1) &= P(|Z| \geq \frac{1.1}{\sqrt{1/6}}) \\ &= 2P(Z < -\frac{1.1}{\sqrt{1/6}}) \approx 0.007 \end{aligned}$$

where $Z \sim \text{Normal}(0, 1)$. As this P-value is less than the significance level, the test rejects the null hypothesis. That is, the test concludes that the slope of 1.1 is far enough from 0 that it demonstrates a true relationship between the X_j input values and the Y_j output values.

EXERCISES

Ex. 10.6.1. Continuing with Example 10.6.1, use Theorem 10.4.3 to devise a hypothesis test for determining whether or not the regression line goes through the origin. That is, determine whether or not $\alpha = 0$ is a plausible assumption.

10.7 ESTIMAING AN UNKNOWN σ^2

In many cases the variance σ^2 of the points around the regression line will be an unknown quantity and so, like α and β , it too will need to be approximated using the $(X_1, Y_1), \dots, (X_n, Y_n)$ data. The following theorem provides an unbiased estimator for σ^2 using the data.

THEOREM 10.7.1. Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be data following the simple linear model with $n > 2$. Let $S^2 = \frac{1}{n-2} \sum_{j=1}^n (Y_j - (a + bX_j))^2$. Then S^2 is an unbiased estimator for σ^2 . (That is, $E[S^2] = \sigma^2$).

Proof - Before looking at $E[S^2]$ in its entirety, we look at three quantities that will be helpful in computing this expected value.

First note,

$$\begin{aligned}
 Var[(Y_j - \bar{Y})] &= Var\left[\frac{nY_j + (Y_1 + Y_2 + \dots + Y_n)}{n}\right] \\
 &= \frac{1}{n^2} \left(Var[(n-1)Y_j + \sum_{i=1, i \neq j}^n Y_i] \right) \\
 &= \frac{1}{n^2} \left([(n-1)^2 \sigma^2 + \sum_{i=1, i \neq j} \sigma^2] \right) \\
 &= \frac{1}{n^2} [(n-1)^2 \sigma^2 + (n-1)\sigma^2] \\
 &= \frac{n-1}{n} \sigma^2
 \end{aligned}$$

and therefore,

$$\begin{aligned}
 \sum_{j=1}^n E[(Y_j - \bar{Y})^2] &= \sum_{j=1}^n Var[Y_j - \bar{Y}] + (E[Y_j - \bar{Y}])^2 \\
 &= \sum_{j=1}^n \frac{n-1}{n} \sigma^2 + ((\alpha + \beta X_j) - (\alpha + \beta \bar{X}))^2 \\
 &= \sum_{j=1}^n \frac{n-1}{n} \sigma^2 + \beta^2 (X_j - \bar{X})^2 \\
 &= (n-1)\sigma^2 + \beta^2 \sum_{j=1}^n (X_j - \bar{X})^2 \\
 &= (n-1)\sigma^2 + \beta^2(n-1)S_X^2. \tag{10.7.1}
 \end{aligned}$$

Next,

$$\begin{aligned}
 \sum_{j=1}^n E[b^2(X_j - \bar{X})^2] &= E[b^2] \sum_{j=1}^n (X_j - \bar{X})^2 \\
 &= (Var[b] + (E[b])^2)((n-1)S_X^2) \\
 &= \left(\frac{\sigma^2}{(n-1)S_X^2} + \beta^2\right)((n-1)S_X^2) \\
 &= \sigma^2 + \beta^2(n-1)S_X^2. \tag{10.7.2}
 \end{aligned}$$

Also,

$$\begin{aligned}
E[bY_j] &= Cov[b, Y_j] + E[b]E[Y_j] \\
&= Cov\left[\sum_{i=1}^n \frac{X_i - \bar{X}}{(n-1)S_X^2} Y_i, Y_j\right] + \beta(\alpha + \beta X_j) \\
&= \sum_{i=1}^n \frac{X_i - \bar{X}}{(n-1)S_X^2} Cov[Y_i, Y_j] + \beta(\alpha + \beta X_j) \\
&= \frac{X_i - \bar{X}}{(n-1)S_X^2} Var[Y_j] + \beta(\alpha + \beta X_j) \\
&= \frac{X_i - \bar{X}}{(n-1)S_X^2} \sigma^2 + \beta(\alpha + \beta X_j)
\end{aligned}$$

from which we may determine that

$$\begin{aligned}
\sum_{j=1}^n E[(Y_j - \bar{Y})b(X_j - \bar{X})] &= \sum_{j=1}^n (X_j - \bar{X})E[Y_j b] - \sum_{j=1}^n (X_j - \bar{X})E[\bar{Y}b] \\
&= \sum_{j=1}^n \sum_{i=1}^n (X_j - \bar{X})\left(\frac{X_i - \bar{X}}{(n-1)S_X^2} \sigma^2 + \beta(\alpha + \beta X_j)\right) - \sum_{j=1}^n (X_j - \bar{X})E[\bar{Y}]E[b] \\
&= \sum_{j=1}^n \frac{(X_j - \bar{X})^2}{(n-1)S_X^2} \sigma^2 + \sum_{j=1}^n (X_j - \bar{X})\beta(\alpha + \beta X_j) - \sum_{j=1}^n (X_j - \bar{X})(\alpha + \beta \bar{X})\beta \\
&= \sigma^2 + \sum_{j=1}^n (X_j - \bar{X})\beta^2(X_j - \bar{X}) \\
&= \sigma^2 + \beta^2(n-1)S_X^2
\end{aligned} \tag{10.7.3}$$

Finally, putting together the results from equations 10.7.1, 10.7.2, and 10.7.3 we find

$$\begin{aligned}
E\left[\sum_{j=1}^n (Y_j - (a + bX_j))^2\right] &= E\left[\sum_{j=1}^n (Y_j - (\bar{Y} + b(X_j - \bar{X})))\right] \\
&= E\left[\sum_{j=1}^n ((Y_j - \bar{Y}) - b(X_j - \bar{X}))^2\right] \\
&= E\left[\sum_{j=1}^n (Y_j - \bar{Y})^2 - 2(Y_j - \bar{Y})b(X_j - \bar{X}) + b^2(X_j - \bar{X})^2\right] \\
&= \sum_{j=1}^n E[(Y_j - \bar{Y})^2] - 2E[(Y_j - \bar{Y})b(X_j - \bar{X})] + E[b^2(X_j - \bar{X})^2] \\
&= ((n-1)\sigma^2 + \beta^2(n-1)S_X^2) - 2(\sigma^2 + \beta^2(n-1)S_X^2) + (\sigma^2 + \beta^2(n-1)S_X^2) \\
&= (n-2)\sigma^2
\end{aligned}$$

Hence $E[S_X^2] = E\left[\frac{1}{n-2} \sum_{j=1}^n (Y_j - (a + bX_j))^2\right] = \sigma^2$ as desired.

4

SUMMARIZING DISCRETE RANDOM VARIABLES

When we first looked at Bernoulli trials in Example 2.1.2 we asked the question “On average how many successes will there be after n trials?” In order to answer this question, a specific definition of “average” must be developed.

To begin, consider how to extend the basic notion of the average of a list of numbers to the situation of equally likely outcomes. For instance, if we want to know what the average roll of a die will be, it makes sense to declare it to be 3.5, the average value of 1, 2, 3, 4, 5, and 6. A motivation for a more general definition of average comes from a rewriting of this calculation.

$$\frac{1+2+3+4+5+6}{6} = 1\left(\frac{1}{6}\right) + 2\left(\frac{1}{6}\right) + 3\left(\frac{1}{6}\right) + 4\left(\frac{1}{6}\right) + 5\left(\frac{1}{6}\right) + 6\left(\frac{1}{6}\right).$$

From the perspective of the right hand side of the equation, the results of all outcomes are added together after being weighted, each according to its probability. In the case of a die, all six outcomes have probability $\frac{1}{6}$.

4.1 EXPECTED VALUE

DEFINITION 4.1.1. Let $X : S \rightarrow T$ be a discrete random variable (so T is countable). Then the expected value (or average) of X is written as $E[X]$ and is given by

$$E[X] = \sum_{t \in T} t \cdot P(X = t)$$

provided that the sum converges absolutely. In this case we say that X has “finite expectation”. If the sum diverges to $\pm\infty$ we say the random variable has infinite expectation. If the sum diverges, but not to infinity, we say the expected value is undefined.

EXAMPLE 4.1.2. In the previous chapter, Example 3.1.4 described a lottery for which a ticket could be worth nothing, or it could be worth either \$20 or \$200. What is the average value of such a ticket?

We calculated the distribution of ticket values as $P(X = 200) = \frac{1}{1000}$, $P(X = 20) = \frac{27}{1000}$, and $P(X = 0) = \frac{972}{1000}$. Applying the definition of expected value results in

$$E[X] = 200\left(\frac{1}{1000}\right) + 20\left(\frac{27}{1000}\right) + 0\left(\frac{972}{1000}\right) = 0.74,$$

so a ticket has an expected value of 56 cents. ■

It is possible to think of a constant as a random variable. If $c \in \mathbb{R}$ then we could define a random variable X with a distribution such that $P(X = c) = 1$. It is a slight abuse of notation, but in this case we will simply write c for both the real number as well as the constant random variable. Such random variables have the obvious expected value.

THEOREM 4.1.3. *Let c be a real number. Then $E[c] = c$.*

Proof - By definition $E[c]$ is a sum over all possible values of c , but in this case that is just a single value, so $E[c] = c \cdot P(c = c) = c \cdot 1 = c$. ■

When the range of X is finite, $E[X]$ always exists since it is a finite sum. When the range of X is infinite there is a possibility that the infinite series will not be absolutely convergent and therefore that $E[X]$ will be infinite or undefined. In fact, when proving theorems about how expected values behave, most of the complications arise from the fact that one must know that an infinite sum converges absolutely in order to rearrange terms within that sum with equality. The next examples explore ways in which expected values may misbehave.

EXAMPLE 4.1.4. Suppose X is a random variable taking values in the range $T = \{2, 4, 8, 16, \dots\}$ such that $P(X = 2^n) = \frac{1}{2^n}$ for all integers $n \geq 1$.

This is the distribution of a random variable since

$$\sum_{n=1}^{\infty} P(X = 2^n) = \sum_{n=1}^{\infty} \frac{1}{2^n} = 1.$$

But note that

$$\sum_{n=1}^{\infty} 2^n \cdot P(X = 2^n) = \sum_{n=1}^{\infty} 2^n \frac{1}{2^n} = \sum_{n=1}^{\infty} 1$$

which diverges to infinity, so this random variable has an infinite expected value. ■

EXAMPLE 4.1.5. Suppose X is a random variable taking values in the range $T = \{-2, 4, -8, 16, \dots\}$ such that $P(X = (-2)^n) = \frac{1}{2^n}$ for all integers $n \geq 1$.

$$\sum_{n=1}^{\infty} (-2)^n \cdot P(X = 2^n) = \sum_{n=1}^{\infty} (-2)^n \frac{1}{2^n} = \sum_{n=1}^{\infty} (-1)^n.$$

This infinite sum diverges (not to $\pm\infty$), so the expected value of this random variable is undefined. ■

The examples above were specifically constructed to produce series which clearly diverged, but in general it can be complicated to check whether an infinite sum is absolutely convergent or not. The next technical lemma provides a condition that is often simpler to check. The convenience of this lemma is that, since $|X|$ is always positive, the terms of the series for $E[|X|]$ may be freely rearranged without changing the value of (or the convergence of) the sum.

LEMMA 4.1.6. *$E[X]$ is a real number if and only if $E[|X|] < \infty$.*

Proof - Let T be the range of X . So $U = \{|t| : t \in T\}$ is the range of $|X|$. By definition

$$E[|X|] = \sum_{u \in U} u \cdot P(|X| = u), \quad \text{while}$$

$$E[X] = \sum_{t \in T} t \cdot P(X = t).$$

To more easily relate these two sums, define $\hat{T} = \{t : |t| \in U\}$. Since every $u \in U$ came from some $t \in T$ the new set \hat{T} contains every element of T . For every $t \in \hat{T}$ for which $t \notin T$, the element is outside of the range of X and so $P(X = t) = 0$ for such elements. Because of this $E[X]$ may be written as

$$E[X] = \sum_{t \in \hat{T}} t \cdot P(X = t)$$

since any additional terms in the series are zero.

Note that for each $u \in U$, the event ($|X| = u$) is equal to $(X = u) \cup (X = -u)$ where each of u and $-u$ is an element of \hat{T} . Therefore,

$$\begin{aligned} u \cdot P(U = u) &= u \cdot (P(X = u) + P(X = -u)) \\ &= u \cdot P(X = u) + u \cdot P(X = -u) \\ &= |u| \cdot P(X = u) + |-u| \cdot P(X = -u) \end{aligned}$$

(When $u = 0$ the quantities $P(|X| = 0)$ and $P(X = 0) + P(X = -0)$ are typically not equal, but the equation is still true since both sides of the equation are zero). Summing over all $u \in U$ then yields

$$\begin{aligned} \sum_{u \in U} u \cdot P(|X| = u) &= \sum_{u \in U} |u| \cdot P(X = u) + |-u| \cdot P(X = -u) \\ &= \sum_{t \in \hat{T}} |t| \cdot P(X = t) \\ &= \sum_{t \in T} |t \cdot P(X = t)|. \end{aligned}$$

Therefore the series describing $E[X]$ is absolutely convergent exactly when $E[|X|] < \infty$. ■

4.1.1 Properties of the Expected Value

We will eventually wish to calculate the expected values of functions of multiple random variables. Of particular interest to statistics is an understanding of expected values of sums and averages of i.i.d. sequences. That understanding will be made easier by first learning something about how expected values behave for simple combinations of variables.

THEOREM 4.1.7. *Suppose that X and Y are discrete random variables, both with finite expected value and both defined on the same sample space S . If a and b are real numbers then*

- (1) $E[aX] = aE[X]$;
- (2) $E[X + Y] = E[X] + E[Y]$; and
- (3) $E[aX + bY] = aE[X] + bE[Y]$.
- (4) If $X \geq 0$ then $E[X] \geq 0$.

Proof of (1) - If $a = 0$ then both sides of the equation are zero, so assume $a \neq 0$. We know that X is a function from S to some range U . So aX is also a random variable and its range is $T = \{au : u \in U\}$.

By definition $E[aX] = \sum_{t \in T} t \cdot P(aX = t)$, but because of how T is defined, adding values indexed by $t \in T$ is equivalent to adding values indexed by $u \in U$ where $t = au$. In other words

$$\begin{aligned} E[aX] &= \sum_{t \in T} t \cdot P(aX = t) \\ &= \sum_{u \in U} au \cdot P(aX = au) \\ &= a \cdot \sum_{u \in U} u \cdot P(X = u) \\ &= aE[X]. \end{aligned}$$

Proof of (2) - We are assuming that X and Y have the same domain, but they typically have different ranges. Suppose $X : S \rightarrow U$ and $Y : S \rightarrow V$. Then the random variable $X + Y$ is also defined on S and takes values in $T = \{u + v : u \in U, v \in V\}$. Therefore, adding values indexed by $t \in T$ is equivalent to adding values indexed by u and v as they range over U and V respectively. So,

$$\begin{aligned} E[X + Y] &= \sum_{t \in T} t \cdot P(X + Y = t) \\ &= \sum_{u \in U, v \in V} (u + v) \cdot P(X = u, Y = v) \\ &= \sum_{u \in U} \sum_{v \in V} (u + v) \cdot P(X = u, Y = v) \\ &= \sum_{u \in U} \sum_{v \in V} u \cdot P(X = u, Y = v) + \sum_{u \in U} \sum_{v \in V} v \cdot P(X = u, Y = v) \\ &= \sum_{u \in U} \sum_{v \in V} u \cdot P(X = u, Y = v) + \sum_{v \in V} \sum_{u \in U} v \cdot P(X = u, Y = v) \end{aligned}$$

where the rearrangement of summation is legitimate since the series converges absolutely. Notice that as u ranges over all of U the sets $(X = u, Y = v)$ partition the set $(Y = v)$ into disjoint pieces based on the value of X . Likewise the event $(X = u)$ is partitioned by $(X = u, Y = v)$ as v ranges over all values of $v \in V$. Therefore, as a disjoint union,

$$(Y = v) = \bigcup_{u \in U} (X = u, Y = v) \quad \text{and} \quad (X = u) = \bigcup_{v \in V} (X = u, Y = v),$$

and so

$$P(Y = v) = \sum_{u \in U} P(X = u, Y = v) \quad \text{and} \quad P(X = u) = \sum_{v \in V} P(X = u, Y = v).$$

From there the proof may be completed, since

$$\begin{aligned} E[X + Y] &= \sum_{u \in U} u \sum_{v \in V} P(X = u, Y = v) + \sum_{v \in V} v \sum_{u \in U} P(X = u, Y = v) \\ &= \sum_{u \in U} u \cdot P(X = u) + \sum_{v \in V} v \cdot P(Y = v) \\ &= E[X] + E[Y]. \end{aligned}$$

Proof of (3) - This is an easy consequence of (1) and (2). From (2) the expected value $E[aX + bY]$ may be rewritten as $E[aX] + E[bY]$. From there, applying (1) shows this is also equal to $aE[X] + bE[Y]$. (Using induction this theorem may be extended to any finite linear combination of random variables, a fact which we leave as an exercise below).

Proof of (4) - We know that X is a function from S to T where $t \in T$ implies that $t \geq 0$. As,

$$E[X] = \sum_{t \in T} t \cdot P(X = t),$$

it follows by definition of series (in the case T is countable) that $E[X] \geq 0$. ■

EXAMPLE 4.1.8. What is the average value of the sum of a pair of dice?

To answer this question by appealing to the definition of expected value would require summing over the eleven possible outcomes $\{2, 3, \dots, 12\}$ and computing the probabilities of each of those outcomes. Theorem 4.1.7 makes things much simpler. We began this section by noting that a single die roll has an expected value of 3.5. The sum of two dice is $X + Y$ where each of X and Y represents the outcome of a single die. So the average value of the sum of a pair of dice is $E[X + Y] = E[X] + E[Y] = 3.5 + 3.5 = 7$. ■

EXAMPLE 4.1.9. Consider a game in which a player might either gain or lose money based on the result. A game is considered “fair” if it is described by a random variable with an expected value of zero. Such a game is fair in the sense that, on average, the player will have no net change in money after playing.

Suppose a particular game is played with one player (the roller) throwing a die. If the die comes up an even number, the roller wins that dollar amount from his opponent. If the die is odd, the roller wins nothing. Obviously the game as stated is not “fair” since the roller cannot lose money and may win something. How much should the roller pay his opponent to play this game in order to make it a fair game?

Let X be the amount of money the rolling player gains by the result on the die. The set of possible outcomes is $T = \{0, 2, 4, 6\}$ and it should be routine at this point to verify that $E[X] = 2$. Let c be the amount of money the roller should pay to play in order to make the game fair. Since X is the amount of money gained by the roll, the net change of money for the roller is $X - c$ after accounting for how much was paid to play. A fair game requires

$$0 = E[X - c] = E[X] - E[c] = 2 - c.$$

So the roller should pay his opponent \$2 to make the game fair. ■

4.1.2 Expected Value of a Product

Theorem 4.1.7 showed that $E[X + Y] = E[X] + E[Y]$. It is natural to ask whether a similar rule exists for the product of variables. While it is not generally the case that the expected value of a product is the product of the expected values, if X and Y happen to be independent, the result is true.

THEOREM 4.1.10. Suppose that X and Y are discrete random variables, both with finite expected value and both defined on the same sample space S . If X and Y are independent, then $E[XY] = E[X]E[Y]$.

Proof - Suppose $X : S \rightarrow U$ and $Y : S \rightarrow V$. Then the random variable XY takes values in $T = \{uv : u \in U, v \in V\}$. So,

$$\begin{aligned} E[XY] &= \sum_{t \in T} t \cdot P(XY = t) \\ &= \sum_{u \in U} \sum_{v \in V} (uv) \cdot P(X = u, Y = v) \\ &= \sum_{u \in U} \sum_{v \in V} (uv) \cdot P(X = u)P(Y = v) \\ &= \sum_{u \in U} u \cdot P(X = u) \sum_{v \in V} v \cdot P(Y = v) \\ &= \left(\sum_{u \in U} u \cdot P(X = u) \right) \left(\sum_{v \in V} v \cdot P(Y = v) \right) \\ &= E[X]E[Y]. \end{aligned}$$

Before showing an example of how this theorem might be used, we provide a demonstration that the result will not typically hold without the assumption of independence.

EXAMPLE 4.1.11. Let $X \sim \text{Uniform}(\{1, 2, 3\})$ and let $Y = 4 - X$. It is easy to verify $Y \sim \text{Uniform}(\{1, 2, 3\})$ as well, but X and Y are certainly dependent. A routine computation shows $E[X] = E[Y] = 2$, and so $E[X]E[Y] = 4$.

However, the random variable XY can only take on two possible values. It may equal 3 (if either $X = 1$ and $Y = 3$ or vice versa) or it may equal 4 (if $X = Y = 2$). So, $P(XY = 3) = \frac{2}{3}$ and $P(XY = 4) = \frac{1}{3}$. Therefore,

$$E[XY] = 3\left(\frac{2}{3}\right) + 4\left(\frac{1}{3}\right) = \frac{10}{3} \neq 4.$$

The conclusion of Theorem 4.1.10 fails since X and Y are dependent. ■

EXAMPLE 4.1.12. Suppose an insurance company assumes that, for a given month, both the number of customer claims X and the average cost per claim Y are independent random variables. Suppose further the company is able to estimate that $E[X] = 100$ and $E[Y] = \$1,250$. How should the company estimate the total cost of all claims that month?

The total cost should be the number of claims times the average cost per claim, or XY . Using Theorem 4.1.10 the expected value of XY is simply the product of the separate expected values.

$$E[XY] = E[X]E[Y] = 100 \cdot \$1,250 = \$125,000.$$

Notice, though, that the assumption of independence played a critical role in this computation. Such an assumption might not be valid for many practical problems. Consider, for example, if a weather event such as a tornado tends to cause both a larger-than-average number of claims and also a larger-than-average value per claim. This could cause the variables X and Y to be dependent and, in such a case, estimating the total cost would not be as simple as taking the product of the separate expected values. ■

4.1.3 Expected Values of Common Distributions

A quick glance at the definition of expected value shows that it only depends on the distribution of the random variable. Therefore one can compute the expected values for the various common distributions we defined in the previous chapter.

EXAMPLE 4.1.13. (Expected Value of a Bernoulli(p))

Let $X \sim \text{Bernoulli}(p)$. So $P(X = 0) = 1 - p$ and $P(X = 1) = p$.

Therefore $E[X] = 0(1 - p) + 1(p) = p$. ■

EXAMPLE 4.1.14. (Expected Value of a Binomial(n,p))

We will show two ways to calculate this expected value – the first is more computationally complicated, but follows from the definition of the binomial distribution directly; the second is simpler, but requires using the relationship between the binomial and Bernoulli random variables. In algebraic terms, if $Y \sim \text{Binomial}(n, p)$ then

$$\begin{aligned} E[Y] &= \sum_{k=0}^n k \cdot P(Y = k) \\ &= \sum_{k=1}^n k \cdot \binom{n}{k} p^k (1-p)^{n-k} \\ &= \sum_{k=1}^n k \cdot \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k} \\ &= np \cdot \sum_{k=1}^n \frac{(n-1)!}{(k-1)!((n-1)-(k-1))!} p^{k-1} (1-p)^{(n-1)-(k-1)} \\ &= np \cdot \sum_{k=1}^n \binom{n-1}{k-1} p^{k-1} (1-p)^{(n-1)-(k-1)} \\ &= np \cdot \sum_{k=0}^{n-1} \binom{n-1}{k} p^k (1-p)^{(n-1)-k} \end{aligned}$$

where the last equality is a shift of variables. But now, by the binomial theorem, the sum $\sum_{k=0}^{n-1} \binom{n-1}{k} p^k (1-p)^{(n-1)-k}$ is equal to 1 and therefore $E[Y] = np$.

Alternatively, recall that the binomial distribution first came about as the total number of successes in n independent Bernoulli trials. Therefore a $\text{Binomial}(n, p)$ distribution results from adding together n independent $\text{Bernoulli}(p)$ random variables. Let X_1, X_2, \dots, X_n be i.i.d. $\text{Bernoulli}(p)$ and let $Y = X_1 + X_2 + \dots + X_n$. Then $Y \sim \text{Binomial}(n, p)$ and

$$\begin{aligned} E[Y] &= E[X_1 + X_2 + \dots + X_n] \\ &= E[X_1] + E[X_2] + \dots + E[X_n] \\ &= p + p + \dots + p = np. \end{aligned}$$

This also provides the answer to part (d) of Example 2.1.2. The expected number of successes in a series of n independent $\text{Bernoulli}(p)$ trials is np . ■

In the next example we will calculate the expected value of a geometric random variable. The computation illustrates a common technique from calculus for simplifying power series by differentiating the sum term-by-term in order to rewrite a complicated series in a simpler way.

EXAMPLE 4.1.15. (Expected Value of a Geometric(p))

If $X \sim \text{Geometric}(p)$ and $0 < p < 1$, then

$$E[X] = \sum_{k=1}^{\infty} k \cdot p(1-p)^{k-1}$$

To evaluate the sum of the series we will need to work the partial sums of the same. For any $n \geq 1$, let

$$\begin{aligned} T_n &= \sum_{k=1}^n kp(1-p)^{k-1} = \sum_{k=1}^n k(1-(1-p))(1-p)^{k-1} \\ &= \sum_{k=1}^n k(1-p)^{k-1} - \sum_{k=1}^n k(1-p)^k \\ &= \sum_{k=1}^n (1-p)^{k-1} - n(1-p)^n = \frac{1-(1-p)^n}{p} - n(1-p)^n. \end{aligned}$$

Using standard results from analysis we know that for $0 < p < 1$,

$$\lim_{n \rightarrow \infty} (1-p)^n = 0 \text{ and } \lim_{n \rightarrow \infty} n(1-p)^n = 0.$$

Therefore $T_n \rightarrow \frac{1}{p}$ as $n \rightarrow \infty$. Hence

$$E[X] = \frac{1}{p}.$$

For instance, suppose we wanted to know on average how many rolls of a die it would take before we observed a 5. Each roll is a Bernoulli trial with a probability $\frac{1}{6}$ of success. The time it takes to observe the first success is distributed as a $Geometric(\frac{1}{6})$ and so has expected value $\frac{1}{1/6} = 6$. On average it should take six rolls before observing this outcome. ■

EXAMPLE 4.1.16. (Expected Value of a Poisson(λ))

We can make a reasonable guess at the expected value of a $Poisson(\lambda)$ random variable by recalling that such a distribution was created to approximate a binomial when n was large and p was small. The parameter $\lambda = np$ remained fixed as we took a limit. Since we showed above that a $Binomial(n, p)$ has an expected value of np , it seems plausible that a $Poisson(\lambda)$ should have an expected value of λ . This is indeed true and it is possible to prove the fact by using the idea that the Poisson random variable is the limit of a sequence of binomial random variables. However, this proof requires an understanding of how limits and expected values interact, a concept that has not yet been introduced in the text. Instead we leave a proof based on a direct algebraic computation as Exercise 4.1.12.

Taking the result as a given, we will illustrate how this expected value might be used for an applied problem. Suppose an insurance company wants to model catastrophic floods using a $Poisson(\lambda)$ random variable. Since floods are rare in any given year, and since the company is considering what might occur over a long span of years, this may be a reasonable assumption.

As its name implies a “50-year flood” is a flood so substantial that it should occur, on average, only once every fifty years. However, this is just an average; it may be possible to have two “50-year floods” in consecutive years, though such an event would be quite rare. Suppose the insurance company wants to know how likely it is that there will be two or more “50-year floods” in the next decade, how should this be calculated?

There is an average of one such flood every fifty years, so by proportional reasoning, in the next ten years there should be an average of 0.2 floods. In other words, the number of floods in the next ten years should a random variable $X \sim Poisson(0.2)$ and we wish to calculate $P(X \geq 2)$.

$$\begin{aligned} P(X \geq 2) &= 1 - P(X = 0) - P(X = 1) \\ &= 1 - e^{-0.2} - e^{-0.2}(0.2) \\ &\approx 0.0002. \end{aligned}$$

So assuming the Poisson random variable is an accurate model, there is only about a 0.02% chance that two or more such disastrous floods would occur in the next decade. ■

For a hypergeometric random variable, we will demonstrate another proof technique common to probability. An expected value may involve a complicated (or infinite) sum which must be computed. However, this sum includes within it the probabilities of each outcome of the random variable, and those probabilities must therefore add to 1. It is sometimes possible to simplify the sum describing the expected value using the fact that a related sum is already known.

EXAMPLE 4.1.17. (Expected Value of a HyperGeo(N, r, m)) Let m and r be positive integers and let N be an integer for which $N > \max\{m, r\}$. Let X be a random variable with $X \sim \text{HyperGeo}(N, r, m)$. To calculate the expected value of X , we begin with two facts. The first is an identity involving combinations. If $n \geq k > 0$ then

$$\begin{aligned} \binom{n}{k} &= \frac{n!}{k!(n-k)!} \\ &= \frac{n}{k} \frac{(n-1)!}{(k-1)!((n-1)-(k-1))!} \\ &= \frac{n}{k} \binom{n-1}{k-1}. \end{aligned}$$

The second comes from the consideration of the probabilities associated with a HyperGeo($N - 1, r - 1, m - 1$) distribution. Specifically, as k ranges over all possible values of such a distribution, we have

$$\sum_k \frac{\binom{r-1}{k} \binom{(N-1)-(r-1)}{(m-1)-k}}{\binom{N-1}{m-1}} = 1$$

since this is the sum over all outcomes of the random variable.

To calculate $E[X]$, let j range over the possible values of X . Recall that the minimum value of j is $\max\{0, m - (N - r)\}$ and the maximum value of j is $\min\{r, m\}$. Now let $k = j - 1$. This means that the maximum value for k is $\min\{r - 1, m - 1\}$. If the minimum value for j was $m - (N - r)$ then the minimum value for k is $m - (N - r) - 1 = ((m - 1) - ((N - 1) - (r - 1)))$. If the minimum value for j was 0 then the minimum value for k is -1 .

The key to the computation is to note that as j ranges over all of the values of X , the values of k cover all possible values of a HyperGeo($N - 1, m - 1, r - 1$) distribution. In fact, the only possible value k may assume that is not in the range of such a distribution is if $k = -1$ as a minimum value. Now,

$$E[X] = \sum_j j \cdot \frac{\binom{r}{j} \binom{N-r}{m-j}}{\binom{N}{m}},$$

and if $j = 0$ is in the range of X , then that term of the sum is zero and it may be deleted without affecting the value. That is equivalent to deleting the $k = -1$ term, so the remaining values of

k exactly describe the range of a HyperGeo($N - 1, m - 1, r - 1$) distribution. From there, the expected value may be calculated as

$$\begin{aligned} E[X] &= \sum_j j \cdot \frac{\binom{r}{j} \binom{N-r}{m-j}}{\binom{N}{m}} \\ &= \sum_j j \cdot \frac{\frac{r}{j} \binom{r-1}{j-1} \binom{(N-1)-(r-1)}{(m-1)-(j-1)}}{\frac{N}{m} \binom{N-1}{m-1}} \\ &= \left(\frac{rm}{N}\right) \cdot \sum_j \frac{\binom{r-1}{j-1} \binom{(N-1)-(r-1)}{(m-1)-(j-1)}}{\binom{N-1}{m-1}} \\ &= \left(\frac{rm}{N}\right) \cdot \sum_k \frac{\binom{r-1}{k} \binom{(N-1)-(r-1)}{(m-1)-k}}{\binom{N-1}{m-1}} \\ &= \left(\frac{rm}{N}\right) \cdot (1) = \frac{rm}{N}. \end{aligned}$$

This nearly completes the goal of calculating the expected values of hypergeometric distributions. The only remaining issues are the cases when $m = 0$ and $r = 0$. Since the hypergeometric distribution was only defined when m and r were non-negative integers, and since the proof above requires the consideration of such a distribution for the values $m - 1$ and $r - 1$, the remaining cases must be handled separately. However, they are fairly easy and yield the same result, a fact we leave it to the reader to verify. ■

4.1.4 Expected Value of $f(X_1, X_2, \dots, X_n)$

As we have seen previously, if X is a random variable and if f is a function defined on the possible outputs of X , then $f(X)$ is a random variable in its own right. The expected value of this new random variable may be computed in the usual way from the distribution of $f(X)$, but it is an extremely useful fact that it may also be computed from the distribution of X itself. The next example and theorems illustrate this fact.

EXAMPLE 4.1.18. Returning to a setting first seen in Example 3.3.1 we will let $X \sim \text{Uniform}(\{-2, -1, 0, 1, 2\})$, and let $f(x) = x^2$. How may $E[f(X)]$ be calculated?

We will demonstrate this in two ways – first by appealing directly to the definition, and then using the distribution of X instead of the distribution of $f(X)$. To use the definition of expected value, recall that $f(X) = X^2$ takes values in $\{0, 1, 4\}$ with the following probabilities: $P(f(X) = 0) = \frac{1}{5}$ while $P(f(X) = 1) = P(f(X) = 4) = \frac{2}{5}$. Therefore,

$$E[f(X)] = 0\left(\frac{1}{5}\right) + 1\left(\frac{2}{5}\right) + 4\left(\frac{2}{5}\right) = 2.$$

However, the values of $f(X)$ are completely determined from the values of X . For instance, the event $(f(X) = 4)$ had a probability of $\frac{2}{5}$ because it was the disjoint union of two other events $(X = 2) \cup (X = -2)$, each of which had probability $\frac{1}{5}$. So the term $4\left(\frac{2}{5}\right)$ in the computation above could equally well have been thought of in two pieces

$$\begin{aligned} 4 \cdot P(f(X) = 4) &= 4 \cdot P((X = 2) \cup (X = -2)) \\ &= 4 \cdot (P(X = 2) + P(X = -2)) \\ &= 4 \cdot P(X = 2) + 4 \cdot P(X = -2) \\ &= 2^2 \cdot P(X = 2) + (-2)^2 \cdot P(X = -2), \end{aligned}$$

where the final expression emphasizes that the outcome of 4 resulted either from 2^2 or $(-2)^2$ depending on the value of X . Following a similar plan for the other values of $f(X)$ allows $E[f(X)]$ to be calculated directly from the probabilities of X as

$$\begin{aligned} E[f(X)] &= (-2)^2 \cdot P(X = -2) + (-1)^2 \cdot P(X = -1) + 0^2 \cdot P(X = 0) \\ &\quad + 1^2 \cdot P(X = 1) + 2^2 \cdot P(X = 2) \\ &= 4\left(\frac{1}{5}\right) + 1\left(\frac{1}{5}\right) + 0\left(\frac{1}{5}\right) + 1\left(\frac{1}{5}\right) + 4\left(\frac{1}{5}\right) \\ &= 2, \end{aligned}$$

which gives the same result as the previous computation. ■

The technique of the example above works for any functions as demonstrated by the next two theorems. We first state and prove a version for functions of a single random variable and then deal with the multivariate case.

THEOREM 4.1.19. *Let $X : S \rightarrow T$ be a discrete random variable and define a function $f : T \rightarrow U$. Then the expected value of $f(X)$ may be computed as*

$$E[f(X)] = \sum_{t \in T} f(t) \cdot P(X = t).$$

Proof - By definition $E[f(X)] = \sum_{u \in U} u \cdot P(f(X) = u)$. However, as in the previous example, the event $(f(X) = u)$ may be partitioned according to the input values of X which cause $f(X)$ to equal u . Recall that $f^{-1}(u)$ describes the set of values in T which, when input into the function f , produce the value u . That is, $f^{-1}(u) = \{t \in T : f(t) = u\}$. Therefore,

$$\begin{aligned} (f(X) = u) &= \bigcup_{t \in f^{-1}(u)} (X = t), \quad \text{and so} \\ P(f(X) = u) &= \sum_{t \in f^{-1}(u)} P(X = t). \end{aligned}$$

Putting this together with the definition of $E[f(X)]$ shows

$$\begin{aligned} E[f(X)] &= \sum_{u \in U} u \cdot P(f(X) = u) \\ &= \sum_{u \in U} u \cdot \sum_{t \in f^{-1}(u)} P(X = t) \\ &= \sum_{u \in U} \sum_{t \in f^{-1}(u)} u \cdot P(X = t) \\ &= \sum_{u \in U} \sum_{t \in f^{-1}(u)} f(t) \cdot P(X = t) \\ &= \sum_{t \in T} f(t) \cdot P(X = t), \end{aligned}$$

where the final step is simply the fact that $T = f^{-1}(U)$ and so summing over the values of $t \in T$ is equivalent to grouping them together in the sets $f^{-1}(u)$ and summing over all values in U that may be achieved by $f(X)$. ■

THEOREM 4.1.20. Let X_1, X_2, \dots, X_n be random variables defined on a common sample space S . The X_j variables may have different ranges, so let $X_j : S \rightarrow T_j$. Let f be a function defined for all possible outputs of the X_j variables. Then

$$E[f(X)] = \sum_{t_1 \in T_1, \dots, t_n \in T_n} f(t_1, \dots, t_n) \cdot P(X_1 = t_1, \dots, X_n = t_n).$$

The proof is nearly the same as for the one-variable case. The only difference is that $f^{-1}(u)$ is now a set of vectors of values (t_1, \dots, t_n) , so that the event $(f(X) = u)$ decomposes into events of the form $(X_1 = t_1, \dots, X_n = t_n)$. However, this change does not interfere with the logic of the proof. We leave the details to the reader.

EXERCISES

Ex. 4.1.1. Let X, Y be discrete random variables. Suppose $X \leq Y$ then show that $E[X] \leq E[Y]$.

Ex. 4.1.2. A lottery is held every day, and on any given day there is a 30% chance that someone will win, with each day independent of every other. Let X denote the random variable describing the number of times in the next five days that the lottery will be won.

- (a) What type of random variable (with what parameter) is X ?
- (b) On average (expected value), how many times in the next five days will the lottery be won?
- (c) When the lottery occurs for each of the next five days, what is the most likely number (mode) of days there will be a winner?
- (d) How likely is it the lottery will be won in either one or two of the next five days?

Ex. 4.1.3. A game show contestant is asked a series of questions. She has a probability of 0.88 of knowing the answer to any given question, independently of every other. Let Y denote the random variable describing the number of questions asked until the contestant does not know the correct answer.

- (a) What type of random variable (with what parameter) is Y ?
- (b) On average (expected value), how many questions will be asked until the first question for which the contestant does not know the answer?
- (c) What is the most likely number of questions (mode) that will be asked until the contestant does not know a correct answer?
- (d) If the contestant is able to answer twelve questions in a row, she will win the grand prize. How likely is it that she will know the answers to all twelve questions?

Ex. 4.1.4. Sonia sends out invitations to eleven of her friends to join her on a hike she's planning. She knows that each of her friends has a 59% chance of deciding to join her independently of each other. Let Z denote the number of friends who join her on the hike.

- (a) What type of random variable (with what parameter) is Z ?
- (b) What is the average (expected value) number of her friends that will join her on the hike?
- (c) What is the most likely number (mode) of her friends that will join her on the hike?

- (d) How do your answers to (b) and (c) change if each friend has only a 41% chance of joining her?

Ex. 4.1.5. A player rolls three dice and earns \$1 for each die that shows a 6. How much should the player pay to make this a fair game?

Ex. 4.1.6. (**The St.Petersburg Paradox**) Suppose a game is played whereby a player begins flipping a fair coin and continues flipping it until it comes up heads. At that time the player wins a 2^n dollars where n is the total number of times he flipped the coin. Show that there is no amount of money the player could pay to make this a fair game. (Hint: See Example 4.1.4).

Ex. 4.1.7. Two different investment strategies have the following probabilities of return on \$10,000.

Strategy A has a 20% chance of returning \$14,000, a 35% chance of returning \$12,000, a 20% chance of returning \$10,000, a 15% chance of returning \$8,000, and a 10% chance of returning only \$6,000.

Strategy B has a 25% chance of returning \$12,000, a 35% chance of returning \$11,000, a 25% chance of returning \$10,000, and a 15% chance of returning \$9,000.

- (a) Which strategy has the larger expected value of return?
- (b) Which strategy is more likely to produce a positive return on investment?
- (c) Is one strategy clearly preferable to the other? Explain your reasoning.

Ex. 4.1.8. Calculate the expected value of a $\text{Uniform}(\{1, 2, \dots, n\})$ random variable by following the steps below.

- (a) Prove the numerical fact that $\sum_{j=1}^n j = \frac{n^2+n}{2}$. (Hint: There are many methods to do this. One uses induction).
- (b) Use (a) to show that if $X \sim \text{Uniform}(\{1, 2, \dots, n\})$, then $E[X] = \frac{n+1}{2}$.

Ex. 4.1.9. Use induction to extend the result of Theorem 4.1.7 by proving the following:

If X_1, X_2, \dots, X_n are random variables with finite expectation all defined on the same sample space S and if a_1, a_2, \dots, a_n are real numbers, then

$$E[a_1X_1 + a_2X_2 + \dots + a_nX_n] = a_1E[X_1] + a_2E[X_2] + \dots + a_nE[X_n].$$

Ex. 4.1.10. Suppose X and Y are random variables for which X has finite expected value and Y has infinite expected value. Prove that $X + Y$ has infinite expected value.

Ex. 4.1.11. Suppose X and Y are random variables. Suppose $E[X] = \infty$ and $E[Y] = -\infty$.

- (a) Provide an example to show that $E[X + Y] = \infty$ is possible.
- (b) Provide an example to show that $E[X + Y] = -\infty$ is possible.
- (c) Provide an example to show that $E[X + Y]$ may have finite expected value.

Ex. 4.1.12. Let $X \sim \text{Poisson}(\lambda)$.

- (a) Write an expression for $E[X]$ as an infinite sum.
- (b) Every non-zero term in your answer to (a) should have a λ in it. Factor this λ out and explain why the remaining sum equals 1. (Hint: One way to do this is through the use of infinite series. Another way is to use the idea from Example 4.1.17).

Ex. 4.1.13. A daily lottery is an event that many people play, but for which the likelihood of any given person winning is very small, making a Poisson approximation appropriate. Suppose a daily lottery has, on average, two winners every five weeks. Estimate the probability that next week there will be more than one winner.

4.2 VARIANCE AND STANDARD DEVIATION

As a single number, the average of a random variable may or may not be a good approximation of the values that variable is likely to produce. For example, let X be defined such that $P(X = 10) = 1$, let Y be defined so that $P(Y = 9) = P(Y = 11) = \frac{1}{2}$, and let Z be defined such that $P(Z = 0) = P(Z = 20) = \frac{1}{2}$. It is easy to check that all three of these random variables have an expected value of 10. However the number 10 exactly describes X , is always off from Y by an absolute value of 1 and is always off from Z by an absolute value of 10.

It is useful to be able to quantify how far away a random variable typically is from its average. Put another way, if we think of the expected value as somehow measuring the “center” of the random variable, we would like to find a way to measure the size of the “spread” of the variable about its center. Quantities useful for this are the variance and standard deviation.

DEFINITION 4.2.1. Let X be a random variable with finite expected value. Then the variance of the random variable is written as $\text{Var}[X]$ and is defined as

$$\text{Var}[X] = E[(X - E[X])^2]$$

The standard deviation of X is written as $SD[X]$ and is defined as

$$SD[X] = \sqrt{\text{Var}[X]}$$

Notice that $\text{Var}[X]$ is the average of the square distance of X from its expected value. So if X has a high probability of being far away from $E[X]$ the variance will tend to be large, while if X is very near $E[X]$ with high probability the variance will tend to be small. In either case the variance is the expected value of a squared quantity, and as such is always non-negative. Therefore $SD[X]$ is defined whenever $\text{Var}[X]$ is defined.

If we were to associate units with the random variable X (say *meters*), then the units of $\text{Var}[X]$ would be *meters*² and the units of $SD[X]$ would be *meters*. We will see that the standard deviation is more meaningful as a measure of the “spread” of a random variable while the variance tends to be a more useful quantity to consider when carrying out complex computations.

Informally we will view the standard deviation as a typical distance from average. So if X is a random variable and we calculate that $E[X] = 12$ and $SD[X] = 3$, we might say, “The variable X will typically take on values that are in or near the range 9 – 15, one standard deviation either side of the average”. A goal of this section is to make that language more precise, but at this point it will help with intuition to understand this informal view.

The variance and standard deviation are described in terms of the expected value. Therefore $\text{Var}[X]$ and $SD[X]$ can only be defined if $E[X]$ exists as a real number. However, it is possible that $\text{Var}[X]$ and $SD[X]$ could be infinite even if $E[X]$ is finite (see Exercises). In practical terms, if X has a finite expected value and infinite standard deviation, it means that the random variable

has a clear average, but is so spread out that any finite number underestimates the typical distance of the random variable from its average.

EXAMPLE 4.2.2. As above, let X be a constant variable with $P(X = 10) = 1$. Let Y be such that $P(Y = 9) = P(Y = 11) = \frac{1}{2}$ and let Z be such that $P(Z = 0) = P(Z = 20) = \frac{1}{2}$.

Since X always equals $E[X]$, the quantity $(X - E[X])^2$ is always zero and we can conclude that $Var[X] = 0$ and $SD[X] = 0$. This makes sense given the view of $SD[X]$ as an estimate of how spread out the variable is. Since X is constant it is not at all spread out and so $SD[X] = 0$.

To calculate $Var[Y]$ we note that $(Y - E[Y])^2$ is always equal to 1. Therefore $Var[Y] = 1$ and $SD[Y] = 1$. Again this reaffirms the informal description of the standard deviation; the typical distance between Y and its average is 1.

Likewise $(Z - E[Z])^2$ is always equal to 100. Therefore $Var[Z] = 100$ and $SD[Z] = 10$. The typical distance between Z and its average is 10. ■

EXAMPLE 4.2.3. What are the variance and standard deviation of a die roll?

Before we carry out the calculation, let us use the informal idea of standard deviation to estimate an answer and help build intuition. We know the average of a die roll is 3.5. The closest a die could possibly be to this average is 0.5 (if it were to roll a 3 or a 4) and the furthest it could possibly be is 2.5 (if it were to roll a 1 or a 6). Therefore the standard deviation, a typical distance from average, should be somewhere between 0.5 and 2.5.

To calculate the quantity exactly, let X represent the roll of a die. By definition, $Var[X] = E[(X - 3.5)^2]$, and the values that $(X - 3.5)^2$ may assume are determined by the six values X may take on.

$$\begin{aligned} Var[X] &= E[(X - 3.5)^2] \\ &= \frac{1}{6}(2.5)^2 + \frac{1}{6}(1.5)^2 + \frac{1}{6}(0.5)^2 + \frac{1}{6}(-0.5)^2 + \frac{1}{6}(-1.5)^2 + \frac{1}{6}(-2.5)^2 \\ &= \frac{35}{12}. \end{aligned}$$

So, $SD[X] = \sqrt{\frac{35}{12}} \approx 1.71$ which is near the midpoint of the range of our estimate above. ■

4.2.1 Properties of Variance and Standard Deviation

THEOREM 4.2.4. Let $a \in \mathbb{R}$ and let X be a random variable with finite variance (and thus, with finite expected value as well). Then,

- (a) $Var[aX] = a^2 \cdot Var[X]$;
- (b) $SD[aX] = |a| \cdot SD[X]$;
- (c) $Var[X + a] = Var[X]$; and
- (d) $SD[X + a] = SD[X]$.

Proof of (a) and (b) - $Var[aX] = E[(aX - E[aX])^2]$. Using known properties of expected value this may be rewritten as

$$\begin{aligned} Var[aX] &= E[(aX - aE[X])^2] \\ &= E[a^2(X - E[X])^2] \\ &= a^2 E[(X - E[X])^2] \\ &= a^2 Var[X]. \end{aligned}$$

That concludes the proof of (a). The result from (b) follows by taking square roots of both sides of this equation.

Proof of (c) and (d) - (See Exercises)

The variance may also be computed using a different (but equivalent) formula if $E[X]$ and $E[X^2]$ are known.

THEOREM 4.2.5. *Let X be a random variable for which $E[X]$ and $E[X^2]$ are both finite. Then*

$$\text{Var}[X] = E[X^2] - (E[X])^2.$$

Proof -

$$\begin{aligned}\text{Var}[X] &= E[(X - E[X])^2] \\ &= E[X^2 - 2XE[X] + (E[X])^2] \\ &= E[X^2] - 2E[XE[X]] + E[(E[X])^2].\end{aligned}$$

But $E[X]$ is a constant, so

$$\begin{aligned}\text{Var}[X] &= E[X^2] - 2E[XE[X]] + E[(E[X])^2] \\ &= E[X^2] - 2E[X]E[X] + (E[X])^2 \\ &= E[X^2] - (E[X])^2.\end{aligned}$$

In statistics we frequently want to consider the sum or average of many random variables. As such it is useful to know how the variance of a sum relates to the variances of each variable separately. Toward that goal we have

THEOREM 4.2.6. *If X and Y are independent random variables, both with finite expectation and finite variance, then*

$$(a) \text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]; \text{ and}$$

$$(b) \text{SD}[X + Y] = \sqrt{(\text{SD}[X])^2 + (\text{SD}[Y])^2}.$$

Proof - Using Theorem 4.2.5,

$$\begin{aligned}\text{Var}[X + Y] &= E[(X + Y)^2] - (E[X + Y])^2 \\ &= E[X^2 + 2XY + Y^2] - ((E[X])^2 + 2E[X]E[Y] + (E[Y])^2) \\ &= E[X^2] + 2E[XY] + E[Y^2] - (E[X])^2 - 2E[X]E[Y] - (E[Y])^2.\end{aligned}$$

But by Theorem 4.1.10, $E[XY] = E[X]E[Y]$ since X and Y are independent. So,

$$\begin{aligned}\text{Var}[X + Y] &= E[X^2] - (E[X])^2 + E[Y^2] - (E[Y])^2 \\ &= \text{Var}[X] + \text{Var}[Y].\end{aligned}$$

Part (b) follows immediately after rewriting the variances in terms of standard deviations and taking square roots. As with expected values, this theorem may be generalized to a sum of any finite number of independent random variables using induction. The proof of that fact is left as Exercise 4.2.11.

EXAMPLE 4.2.7. What is the standard deviation of the sum of two dice?

We previously found that if X represents one die, then $\text{Var}[X] = \frac{35}{12}$. If X and Y are two independent dice, then $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] = \frac{35}{12} + \frac{35}{12} = \frac{35}{6}$. Therefore $\text{SD}[X + Y] = \sqrt{\frac{35}{6}} \approx 2.42$.

4.2.2 Variances of Common Distributions

As with expected value, the variances of the common discrete random variables can be calculated from their corresponding distributions.

EXAMPLE 4.2.8. (Variance of a Bernoulli(p))

Let $X \sim \text{Bernoulli}(p)$. We have already calculated that $E[X] = p$. Since X only takes on the values 0 or 1 it is always true that $X^2 = X$. Therefore $E[X^2] = E[X] = p$.

So, $\text{Var}[X] = E[X^2] - (E[X])^2 = p - p^2 = p(1 - p)$. ■

EXAMPLE 4.2.9. (Variance of a Binomial(n,p))

We will calculate the variance of a binomial using the fact that it may be viewed as the sum of n independent Bernoulli random variables. A strictly algebraic computation is also possible (see Exercises).

Let X_1, X_2, \dots, X_n be independent $\text{Bernoulli}(p)$ random variables. Therefore, if $Y = X_1 + X_2 + \dots + X_n$ then $Y \sim \text{Binomial}(n, p)$ and

$$\begin{aligned}\text{Var}[Y] &= \text{Var}[X_1 + X_2 + \dots + X_n] \\ &= \text{Var}[X_1] + \text{Var}[X_2] + \dots + \text{Var}[X_n] \\ &= p(1 - p) + p(1 - p) + \dots + p(1 - p) \\ &= np(1 - p).\end{aligned}$$

For an application of this computation we return to the idea of sampling from a population where some members of the population have a certain characteristic and others do not. The goal is to provide an estimate of the number of people in the sample that have the characteristic. For this example, suppose we were to randomly select 100 people from a large city in which 20% of the population works in a service industry. How many of the 100 people from our sample should we expect to be service industry workers?

If the sampling is done without replacement (so we cannot pick the same person twice), then strictly speaking the desired number would be described by a hypergeometric random variable. However, we have also seen that there is little difference between the binomial and hypergeometric distributions when the size of the sample is small relative to the size of the population. So since the sample is only 100 people from a “large city”, we will assume this situation is modeled by a binomial random variable. Specifically, since 20% of the population consists of service workers, we will assume $X \sim \text{Binomial}(100, 0.2)$.

The simplest way to answer to the question of how many service industry workers to expect within the sample is to compute the expected value of X . In this case $E[X] = 100(0.2) = 20$, so we should expect around 20 of the 100 people in the sample to be service workers. However, this is an incomplete answer to the question since it only provides an average value; the actual number of service workers in the sample is probably not going to be exactly 20, it's only likely to be around 20 on average. A more complete answer to the question would give an estimate as to how far away from 20 the actual value is likely to be. But this is precisely what the standard deviation describes – an estimate of the likely difference between the actual result of the random variable and its expected value.

In this case $\text{Var}[X] = 100(0.2)(0.8) = 16$ and so $\text{SD}[X] = \sqrt{16} = 4$. This means that the actual number of service industry workers in the sample will typically be about 4 or so away from the expected value of 20, so a more complete answer to the question would be “The sample is likely to have around 16 – 24 service workers in it”. That is not to say that the actual number of service workers is guaranteed to fall in the that range, but the range provides a sort of likely

error associated with the estimate of 20. Results in the $16 - 24$ range should be considered fairly common. Results far outside that range, while possible, should be considered fairly unusual. ■

Recall in Example 4.1.17 we calculated $E[X]$ using a technique in which the sum describing $E[X]$ was computed based on another sum which only involved the distribution of X directly. This second sum equalled 1 since it simply added up the probabilities that X assumed each of its possible values. In a similar fashion, it is sometimes possible to calculate a sum describing $E[X^2]$ in terms of a sum for $E[X]$ which is already known. From that point, Theorem 4.2.5 may be used to calculate the variance and standard deviation of X . This technique will be illustrated in the next example in which we calculate the spread associated with a geometric random variable.

EXAMPLE 4.2.10. (Variance of a Geometric(p))

Let $0 < p < 1$. $X \sim \text{Geometric}(p)$ for which we know $E[X] = \frac{1}{p}$. Then,

$$E[X^2] = \sum_{k=1}^{\infty} k^2 p(1-p)^{k-1}$$

To evaluate the sum of the series we will need to work the partial sums of the same. For any $n \geq 1$, let

$$\begin{aligned} S_n &= \sum_{k=1}^n k^2 p(1-p)^{k-1} = \sum_{k=1}^n k^2 (1-(1-p))(1-p)^{k-1} \\ &= \sum_{k=1}^n k^2 (1-p)^{k-1} - \sum_{k=1}^n k^2 (1-p)^k \\ &= 1 + \sum_{k=2}^n (2k-1)(1-p)^{k-1} - n^2 (1-p)^n \\ &= 1 - \sum_{k=2}^n (1-p)^{k-1} + 2 \sum_{k=2}^n k(1-p)^{k-1} - n^2 (1-p)^n \\ &= 2 - \sum_{k=1}^n (1-p)^{k-1} + 2(-1 + \sum_{k=1}^n k(1-p)^{k-1}) - n^2 (1-p)^n \\ &= -\frac{1-(1-p)^n}{p} + \frac{2}{p} \sum_{k=1}^n kp(1-p)^{k-1} - n^2 (1-p)^n \end{aligned}$$

Using standard results from analysis and result from Example 4.1.15 we know that for $0 < p < 1$,

$$\lim_{n \rightarrow \infty} \sum_{k=1}^n kp(1-p)^{k-1} = \frac{1}{p}, \lim_{n \rightarrow \infty} (1-p)^n = 0, \text{ and } \lim_{n \rightarrow \infty} n^2 (1-p)^n = 0.$$

Therefore $S_n \rightarrow -\frac{1}{p} + \frac{2}{p^2}$ as $n \rightarrow \infty$. Hence

$$E[X^2] = -\frac{1}{p} + \frac{2}{p^2}.$$

Using Theorem 4.2.5 the variance may then be calculated as

$$\begin{aligned} \text{Var}[X] &= E[X^2] - (E[X])^2 \\ &= \frac{2}{p^2} - \frac{1}{p} - \left(\frac{1}{p}\right)^2 \\ &= \frac{1}{p^2} - \frac{1}{p} \end{aligned}$$

■

A similar technique may be used for calculating the variance of a Poisson random variable, a fact which is left as an exercise. We finish this subsection with a computation of the variance of a hypergeometric distribution using an idea similar to how we calculated its expected value in Example 4.1.17.

EXAMPLE 4.2.11. Let m and r be positive integers and let N be an integer with $N > \max\{m, r\}$ and let $X \sim \text{HyperGeo}(N, r, m)$. To calculate $E[X^2]$, as j ranges over the values of X ,

$$\begin{aligned} E[X^2] &= \sum_j j^2 \cdot \frac{\binom{r}{j} \binom{N-r}{m-j}}{\binom{N}{m}} \\ &= \sum_j j^2 \cdot \frac{\frac{r}{j} \binom{r-1}{j-1} \binom{(N-1)-(r-1)}{(m-1)-(j-1)}}{\frac{N}{m} \binom{N-1}{m-1}} \\ &= \left(\frac{rm}{N}\right) \sum_j j \cdot \frac{\binom{r-1}{j-1} \binom{(N-1)-(r-1)}{(m-1)-(j-1)}}{\binom{N-1}{m-1}} \\ &= \left(\frac{rm}{N}\right) \cdot \sum_k (k+1) \frac{\binom{r-1}{k} \binom{(N-1)-(r-1)}{(m-1)-k}}{\binom{N-1}{m-1}} \end{aligned}$$

where k ranges over the values of $Y \sim \text{HyperGeo}(N-1, r-1, m-1)$. Therefore,

$$\begin{aligned} E[X^2] &= \left(\frac{rm}{N}\right) E[Y+1] \\ &= \left(\frac{rm}{N}\right) (E[Y] + 1) \\ &= \left(\frac{rm}{N}\right) \left(\frac{(r-1)(m-1)}{(N-1)} + 1\right). \end{aligned}$$

Now the variance may be easily computed as

$$\begin{aligned} \text{Var}[X] &= E[X^2] - (E[X])^2 \\ &= \left(\frac{rm}{N}\right) \left(\frac{(r-1)(m-1)}{(N-1)} + 1\right) - \left(\frac{rm}{N}\right)^2 \\ &= \frac{N^2 rm - N rm^2 - Nr^2 m + r^2 m^2}{N^2(N-1)}. \end{aligned}$$

As with the computation of expected value, the cases of $m = 0$ and $r = 0$ must be handled separately, but yield the same result. ■

4.2.3 Standardized Variables

Many random variables may be rescaled into a standard format by shifting them so that they have an average of zero and then rescaling them so that they have a variance (and standard deviation) of one. We introduce this idea now, though its chief importance will not be realized until later.

DEFINITION 4.2.12. A standardized random variable X is one for which

$$E[X] = 0 \quad \text{and} \quad \text{Var}[X] = 1.$$

THEOREM 4.2.13. Let X be a discrete random variable with finite expected value and finite, non-zero variance. Then $Z = \frac{X - E[X]}{\text{SD}[X]}$ is a standardized random variable.

Proof - The expected value of Z is

$$\begin{aligned} E[Z] &= E\left[\frac{X - E[X]}{\text{SD}[X]}\right] \\ &= \frac{E[X - E[X]]}{\text{SD}[X]} \\ &= \frac{E[X] - E[X]}{\text{SD}[X]} = 0 \end{aligned}$$

while the variance of Z is

$$\begin{aligned} \text{Var}[Z] &= \text{Var}\left[\frac{X - E[X]}{\text{SD}[X]}\right] \\ &= \frac{\text{Var}[X - E[X]]}{(\text{SD}[X])^2} \\ &= \frac{\text{Var}[X]}{\text{Var}[X]} = 1. \end{aligned}$$

For easy reference we finish off this section by providing a chart of values associated with common discrete distributions.

Distribution	Expected Value	Variance
Bernoulli(p)	p	$p(1-p)$
Binomial(n, p)	np	$np(1-p)$
Geometric(p)	$\frac{1}{p}$	$\frac{1-p}{p^2}$
HyperGeo(N, r, m)	$\frac{rm}{N}$	$\frac{N^2 rm - Nrm^2 - Nr^2 m + r^2 m^2}{N^2(N-1)}$
Poisson(λ)	λ	λ
Uniform($\{1, 2, \dots, n\}$)	$\frac{n+1}{2}$	$\frac{n^2-1}{12}$

EXERCISES

Ex. 4.2.1. A random variable X has a probability mass function given by

$$P(X = 0) = 0.2, P(X = 1) = 0.5, P(X = 2) = 0.2, \text{ and } P(X = 3) = 0.1.$$

Calculate the expected value and standard deviation of this random variable. What is the probability this random variable will produce a result more than one standard deviation from its expected value?

Ex. 4.2.2. Answer the following questions about flips of a fair coin.

- (a) Calculate the standard deviation of the number of heads that show up in 100 flips of a fair coin.
- (b) Show that if the number of coins is quadrupled (to 400) the standard deviation only doubles.

Ex. 4.2.3. Suppose we begin rolling a die, and let X be the number of rolls needed before we see the first 3.

- (a) Show that $E[X] = 6$.
- (b) Calculate $SD[X]$.
- (c) Viewing $SD[X]$ as a typical distance of X from its expected value, would it seem unusual to roll the die more than nine times before seeing a 3?
- (d) Calculate the actual probability $P(X > 9)$.
- (e) Calculate the probability X produces a result within one standard deviation of its expected value.

Ex. 4.2.4. A key issue in statistical sampling is the determination of how much a sample is likely to differ from the population it came from. This exercise explores some of these ideas.

- (a) Suppose a large city is exactly 50% women and 50% men and suppose we randomly select 60 people from this city as part of a sample. Let X be the number of women in the sample. What are the expected value and standard deviation of X ? Given these values, would it seem unusual if fewer than 45% of the individuals in the sample were women?
- (b) Repeat part (a), but now assume that the sample consists of 600 people.

Ex. 4.2.5. Calculate the variance and standard deviation of the value of the lottery ticket from Example 3.1.4.

Ex. 4.2.6. Prove parts (c) and (d) of Theorem 4.2.4.

Ex. 4.2.7. Let $X \sim \text{Binomial}(n, p)$. Show that for $0 < p < 1$, this random variable has the largest standard deviation when $p = \frac{1}{2}$.

Ex. 4.2.8. Follow the steps below to calculate the variance of a random variable with a $\text{Uniform}(\{1, 2, \dots, n\})$ distribution.

- (a) Prove that $\sum_{k=1}^n k^2 = \frac{n(n+1)(2n+1)}{6}$. (Induction is one way to do this).
- (b) Let $X \sim \text{Uniform}(\{1, 2, \dots, n\})$. Use (a) to calculate $E[X^2]$.
- (c) Use (b) and the fact that $E[X] = \frac{n+1}{2}$ to calculate $Var[X]$.

Ex. 4.2.9. This exercise provides an example of a random variable with finite expected value, but infinite variance. Let X be a random variable for which $P(X = \frac{2^n}{n(n+1)}) = \frac{1}{2^n}$ for all integers $n \geq 1$.

- (a) Prove that X is a well-defined variable by showing $\sum_{n=1}^{\infty} P(X = \frac{2^n}{n(n+1)}) = 1$.
- (b) Prove that $E[X] = 1$.

- (c) Prove that $\text{Var}[X]$ is infinite.

Ex. 4.2.10. Recall that the hypergeometric distribution was first developed to answer questions about sampling without replacement. With that in mind, answer the following questions using the chart of expected values and variances.

- Use the formula in the chart to calculate the variance of a hypergeometric distribution if $m = 0$. Explain this result in the context of what it means in terms of sampling.
- Use the formula in the chart to calculate the variance of a hypergeometric distribution if $r = 0$. Explain this result in the context of what it means in terms of sampling.
- Though we only defined a hypergeometric distribution if $N > \max\{r, m\}$, the definition could be extended to $N = \max\{r, m\}$. Use the chart to calculate the variance of a hypergeometric distribution if $N = m$. Explain this result in the context of what it means in terms of sampling without replacement.

Ex. 4.2.11. Prove the following facts about independent random variables.

- (a) Use Theorem 4.2.6 and induction to prove that if X_1, X_2, \dots, X_n are independent, then

$$\text{Var}[X_1 + \dots + X_n] = \text{Var}[X_1] + \dots + \text{Var}[X_n].$$

- (b) Suppose X_1, X_2, \dots, X_n are i.i.d. Prove that

$$E[X_1 + \dots + X_n] = n \cdot E[X_1] \quad \text{and} \quad SD[X_1 + \dots + X_n] = \sqrt{n} \cdot SD[X_1].$$

- (c) Suppose X_1, X_2, \dots, X_n are mutually independent standardized random variables (not necessarily identically distributed). Let

$$Y = \frac{X_1 + X_2 + \dots + X_n}{\sqrt{n}}.$$

Prove that Y is a standardized random variable.

Ex. 4.2.12. Let X be a discrete random variable which takes on only non-negative values. Show that if $E[X] = 0$ then $P(X = 0) = 1$.

Ex. 4.2.13. Suppose X is a discrete random variable with finite variance (and thus finite expected value as well) and suppose there are two different numbers $a, b \in \mathbb{R}$ for which $P(X = a)$ and $P(X = b)$ are both positive. Prove that $\text{Var}[X] > 0$.

Ex. 4.2.14. Let X be a discrete random variable with finite variance (and thus finite expected value as well).

- (a) Prove that $E[X^2] \geq (E[X])^2$.

- (b) Suppose there are two different numbers $a, b \in \mathbb{R}$ for which $P(X = a)$ and $P(X = b)$ are both positive. Prove that $E[X^2] > (E[X])^2$.

Ex. 4.2.15. Let $X \sim \text{Binomial}(n, p)$ for $n > 1$ and $0 < p < 1$. Using the steps below, provide an algebraic proof of the fact that $\text{Var}[X] = np(1 - p)$ without appealing to the fact that such a variable is the sum of Bernoulli trials.

- (a) Begin by writing an expression for $E[X^2]$ in summation form.

- (b) Use (a) to show that $E[X^2] = np \cdot \sum_{k=0}^{n-1} (k+1) \binom{n-1}{k} p^k (1-p)^{(n-1)-k}$.

- (c) Use (b) to explain why $E[X^2] = np \cdot E[Y + 1]$ where $Y \sim \text{Binomial}(n-1, p)$.

- (d) Use (c) together with Theorem 4.2.5 to prove that $\text{Var}[X] = np(1 - p)$.

4.3 STANDARD UNITS

When there is no confusion about what random variable is being discussed, it is usual to use the Greek letter μ in place of $E[X]$ and σ in place of $SD[X]$. When more than one variable is involved the same letters can be used with subscripts (μ_X and σ_X) to indicate which variable is being described.

In statistics one frequently measures results in terms of “standard units” – the number of standard deviations a result is from its expected value. For instance if $\mu = 12$ and $\sigma = 5$, then a result of $X = 20$ would be 1.6 standard units because $20 = \mu + 1.6\sigma$. That is, 20 is 1.6 standard deviations above expected value. Similarly a result of $X = 10$ would be -0.4 standard units because $10 = \mu - 0.4\sigma$.

Since the standard deviation measures a typical distance from average, results that are within one standard deviation from average (between -1 and $+1$ standard units) will tend to be fairly common, while results that are more than two standard deviations from average (less than -2 or greater than $+2$ in standard units) will usually be relatively rare. The likelihoods of some such events will be calculated in the next two examples. Notice that the event ($|X - \mu| \leq k\sigma$) describes those outcomes of X that are within k standard deviations from average.

EXAMPLE 4.3.1. Let Y represent the sum of two dice. How likely is it that Y will be within one standard deviation of its average? How likely is it that Y will be more than two standard deviations from its average?

We can use our previous calculations that $\mu = 7$ and $\sigma = \sqrt{\frac{35}{6}} \approx 2.42$. The achievable values that are within one standard deviation of average are 5, 6, 7, 8, and 9. So the probability that the sum of two dice will be within one standard deviation of average is

$$\begin{aligned} P(|Y - \mu| \leq \sigma) &= P(Y \in \{5, 6, 7, 8, 9\}) \\ &= \frac{4}{36} + \frac{5}{36} + \frac{6}{36} + \frac{5}{36} + \frac{4}{36} \\ &= \frac{2}{3}. \end{aligned}$$

There is about a 66.7% chance that a pair of dice will fall within one standard deviation of their expected value.

Two standard deviations is $2\sqrt{\frac{35}{6}} \approx 4.83$. Only the results 2 and 12 further than this distance from the expected value, so the probability that X will be more than two standard deviations from average is

$$\begin{aligned} P(|Y - \mu| > 2\sigma) &= P(Y \in \{2, 12\}) \\ &= \frac{2}{36} \approx 0.056. \end{aligned}$$

There is only about a 5.6% chance that a pair of dice will be more than two standard deviations from expected value. ■

EXAMPLE 4.3.2. If $X \sim Uniform\{(1, 2, \dots, 100)\}$, what is the probability that X will be within one standard deviation of expected value? What is the probability it will be more than two standard deviations from expected value?

Again, based on earlier calculations we know that $\mu = \frac{101}{2} = 50.5$ and that $\sigma = \sqrt{\frac{9999}{12}} \approx 28.9$. Of the possible values that X can achieve, only the numbers $22, 23, \dots, 79$ fall within one standard deviation of average. So the desired probability is

$$\begin{aligned} P(|X - \mu| \leq \sigma) &= P(X \in \{22, 23, \dots, 79\}) \\ &= \frac{58}{100}. \end{aligned}$$

There is a 58% chance that this random variable will be within one standard deviation of expected value.

Similarly we can calculate that two standard deviations is $2\sqrt{\frac{9999}{12}} \approx 57.7$. Since $\mu = 50.5$ and since the minimal and maximal values of X are 1 and 100 respectively, results that are more than two or more standard deviations from average cannot happen at all for this random variable. In other words $P(|X - \mu| > 2\sigma) = 0$. ■

4.3.1 Markov and Chebyshev Inequalities

The examples of the previous section show that the exact probabilities a random variable will fall within a certain number of standard deviations of its expected value depend on the distribution of the random variable. However, there are some general results that apply to all random variables. To prove these results we will need to investigate some inequalities.

THEOREM 4.3.3. (Markov's Inequality) *Let X be a discrete random variable which takes on only non-negative values and suppose that X has a finite expected value. Then for any $c > 0$,*

$$P(X \geq c) \leq \frac{\mu}{c}.$$

Proof - Let T be the range of X , so T is a countable subset of the positive real numbers. By dividing T into those numbers smaller than c and those numbers that are at least as large as c we have

$$\begin{aligned} \mu &= \sum_{t \in T} t \cdot P(X = t) \\ &= \sum_{t \in T, t < c} t \cdot P(X = t) + \sum_{t \in T, t \geq c} t \cdot P(X = t). \end{aligned}$$

The first sum must be non-negative, since we assumed that T consisted of only non-negative numbers, so we only make the quantity smaller by deleting it. Likewise, for each term in the second sum, $t \geq c$ so we only make the quantity smaller by replacing t by c . This gives us

$$\begin{aligned} \mu &= \sum_{t \in T, t < c} t \cdot P(X = t) + \sum_{t \in T, t \geq c} t \cdot P(X = t) \\ &\geq \sum_{t \in T, t \geq c} c \cdot P(X = t) \\ &= c \cdot \sum_{t \in T, t \geq c} P(X = t). \end{aligned}$$

The events $(X = t)$ indexed over all values $t \in T$ for which $t \geq c$ are a countable collection of disjoint sets whose union is $(X \geq c)$. So,

$$\begin{aligned} \mu &\geq c \cdot \sum_{t \in T, t \geq c} P(X = t) \\ &= cP(X \geq c). \end{aligned}$$

Dividing by c gives the desired result.

Markov's theorem can be useful in its own right for producing an upper bound on the likelihood of certain events, but for now we will use it simply as a lemma to prove our next result.

THEOREM 4.3.4. (Chebychev's Inequality) *Let X be a discrete random variable with finite, non-zero variance. Then for any $k > 0$,*

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

Proof - The event $(|X - \mu| \geq k\sigma)$ is the same as the event $((X - \mu)^2 \geq k^2\sigma^2)$. The random variable $(X - \mu)^2$ is certainly non-negative and its expected value is the variance of X which we have assumed to be finite. Therefore we may apply Markov's inequality to $(X - \mu)^2$ to get

$$\begin{aligned} P(|X - \mu| \geq k\sigma) &= P((X - \mu)^2 \geq k^2\sigma^2) \\ &\leq \frac{E[(X - \mu)^2]}{k^2\sigma^2} \\ &= \frac{\text{Var}[X]}{k^2\sigma^2} \\ &= \frac{\sigma^2}{k^2\sigma^2} \\ &= \frac{1}{k^2}. \end{aligned}$$

Though the theorem is true for all $k > 0$, it doesn't give any useful information unless $k > 1$.

EXAMPLE 4.3.5. Let X be a discrete random variable. Find an upper bound on the likelihood that X will be more than two standard deviations from its expected value.

For the question to make sense we need to assume that X has finite variance to begin with. In which case we may apply Chebychev's inequality with $k = 2$ to find that

$$P(|X - \mu| > 2\sigma) \leq P(|X - \mu| \geq 2\sigma) \leq \frac{1}{4}.$$

There is at most a 25% chance that a random variable will be more than two standard deviations from its expected value. ■

EXERCISES

Ex. 4.3.1. Let $X \sim \text{Binomial}(10, \frac{1}{2})$.

- (a) Calculate μ and σ .
- (b) Calculate $P(|X - \mu| \leq \sigma)$, the probability that X will be within one standard deviation of average. Approximate your answer to the nearest tenth of a percent.
- (c) Calculate $P(|X - \mu| > 2\sigma)$, the probability that X will be more than two standard deviations from average. Approximate your answer to the nearest tenth of a percent.

Ex. 4.3.2. Let $X \sim \text{Geometric}(\frac{1}{4})$.

- (a) Calculate μ and σ .
- (b) Calculate $P(|X - \mu| \leq \sigma)$, the probability that X will be within one standard deviation of average. Approximate your answer to the nearest tenth of a percent.

- (c) Calculate $P(|X - \mu| > 2\sigma)$, the probability that X will be more than two standard deviations from average. Approximate your answer to the nearest tenth of a percent.

Ex. 4.3.3. Let $X \sim Poisson(3)$.

- (a) Calculate μ and σ .
- (b) Calculate $P(|X - \mu| \leq \sigma)$, the probability that X will be within one standard deviation of average. Approximate your answer to the nearest tenth of a percent.
- (c) Calculate $P(|X - \mu| > 2\sigma)$, the probability that X will be more than two standard deviations from average. Approximate your answer to the nearest tenth of a percent.

Ex. 4.3.4. Let $X \sim Binomial(n, \frac{1}{2})$. Determine the smallest value of n for which $P(|X - \mu| > 4\sigma) > 0$. That is, what is the smallest n for which there is a positive probability that X will be more than four standard deviations from average.

Ex. 4.3.5. For $k \geq 1$ there are distributions for which Chebychev's inequality is an equality.

- (a) Let X be a random variable with probability mass function $P(X = 1) = P(X = -1) = \frac{1}{2}$. Prove that Chebychev's inequality is an equality for this random variable when $k = 1$.
- (b) Let X be a random variable with probability mass function $P(X = 1) = P(X = -1) = p$ and $P(X = 0) = 1 - 2p$. For any given value of $k > 1$, show that it is possible to select a value of p for which Chebychev's inequality is an equality when applied to this random variable.

Ex. 4.3.6. Let X be a discrete random variable with finite expected value μ and finite variance σ^2 .

- (a) Explain why $P(|X - \mu| > \sigma) = P((X - \mu)^2 > \sigma^2)$.
- (b) Let T be the range of the random variable $(X - \mu)^2$.
Explain why $\sum_{t \in T} P((X - \mu)^2 = t) = 1$.
- (c) Explain why $Var[X] = \sum_{t \in T} t \cdot P((X - \mu)^2 = t)$.
- (d) Prove that if $P(|X - \mu| > \sigma) = 1$, then
$$Var[X] > \sum_{t \in T} \sigma^2 \cdot P((X - \mu)^2 = t).$$
 (Hint: Use (a) to explain why replacing t by σ^2 in the sum from (c) will only make the quantity smaller).
- (e) Use parts (b) and (d) to derive a contradiction. Note that this proves that the assumption that was made in part (d), namely that $P(|X - \mu| > \sigma) = 1$, cannot be true for any discrete random variable where μ and σ are finite quantities. In other words, no random variable can produce only values that are more than one standard deviation from average.

Ex. 4.3.7. Let X be a discrete random variable with finite expected value and finite variance.

- (a) Prove $P(|X - \mu| \geq \sigma) = 1 \iff P(|X - \mu| = \sigma) = 1$. (A random variable that assumes only values one or more standard deviations from average must only produce values that are exactly one standard deviation from average).
- (b) Prove that if $P(|X - \mu| > \sigma) > 0$ then $P(|X - \mu| < \sigma) > 0$. (If a random variable is able to produce values more than one standard deviation from average, it must also be able to produce values that are less than one standard deviation from average).

4.4 CONDITIONAL EXPECTATION AND CONDITIONAL VARIANCE

In previous chapters we saw that information that a particular event had occurred could substantially change the probability associated with another event. That realization led us to the notion of conditional probability. It is also reasonable to ask how such information might affect the expected value or variance of a random variable.

DEFINITION 4.4.1. Let $X : S \rightarrow T$ be a discrete random variable and let $A \subset S$ be an event for which $P(A) > 0$. The “conditional expected value” is defined from conditional probabilities in the same way the (ordinary) expected value is defined from (ordinary) probabilities. Likewise the “conditional variance” is described in terms of the conditional expected value in the same way the (ordinary) variance is described in terms of the (ordinary) expected value. Specifically, the “conditional expected value” of X given A is

$$E[X|A] = \sum_{t \in T} t \cdot P(X = t|A),$$

and the “conditional variance” of X given A is

$$\text{Var}[X|A] = E[(X - E[X|A])^2|A].$$

EXAMPLE 4.4.2. A die is rolled. What are the expected value and variance of the result given that the roll was even?

Let X be the die roll. Then $X \sim \text{Uniform}(\{1, 2, 3, 4, 5, 6\})$, but conditioned on the event A that the roll was even, this changes so that

$$P(X = 1|A) = P(X = 3|A) = P(X = 5|A) = 0 \quad \text{while}$$

$$P(X = 2|A) = P(X = 4|A) = P(X = 6|A) = \frac{1}{3}.$$

Therefore,

$$E[X|A] = 2\left(\frac{1}{3}\right) + 4\left(\frac{1}{3}\right) + 6\left(\frac{1}{3}\right) = 4.$$

Note that the (unconditioned) expected value of a die roll is $E[X] = 3.5$, so the knowledge of event A slightly increases the expected value of the die roll.

The conditional variance is

$$\text{Var}[X|A] = (2 - 4)^2\left(\frac{1}{3}\right) + (4 - 4)^2\left(\frac{1}{3}\right) + (6 - 4)^2\left(\frac{1}{3}\right) = \frac{8}{3}.$$

This result is slightly less than $\frac{35}{12}$, the (unconditional) variance of a die roll. This means that knowledge of event A slightly decreased the typical spread of the die roll results. ■

In many cases the event A on which an expected value is conditioned will be described in terms of another random variable. For instance $E[X|Y = y]$ is the conditional expectation of X given that variable Y has taken on the value y .

EXAMPLE 4.4.3. Cards are drawn from an ordinary deck of 52, one at a time, randomly and with replacement. Let X and Y denote the number of draws until the first ace and first king are drawn,

respectively. We are interested in say, $E[X|Y = 3]$. When $Y = 3$ an ace was seen on draw 3, but not on draws 1 or 2. Hence

$$P(\text{king on draw } n|Y = 3) = \begin{cases} \frac{4}{48} & \text{if } n = 1 \text{ or } 2 \\ 0 & \text{if } n = 3 \\ \frac{4}{52} & \text{if } n > 3 \end{cases}$$

so that

$$P(X = n|Y = 3) = \begin{cases} \left(\frac{44}{48}\right)^{n-1} \frac{4}{48} & \text{if } n = 1 \text{ or } 2 \\ 0 & \text{if } n = 3 \\ \left(\frac{44}{48}\right)^2 \left(\frac{48}{52}\right)^{n-4} \frac{4}{52} & \text{if } n > 3 \end{cases}$$

For example, when $n > 3$, in order to have $X = n$ a non-king must have been seen on draws 1 and 2 (each with probability $\frac{44}{48}$), a non-king must have resulted on draw 3 (which is automatic, since an ace was drawn), a non-king must have been seen on each of draws 4 through $n - 1$ (each with probability $\frac{48}{52}$), and finally a king was produced on draw n (with probability $\frac{4}{52}$). Hence,

$$\begin{aligned} E[X|Y = 3] &= \sum_{n=1}^2 n \left(\frac{44}{48}\right)^{n-1} \frac{4}{48} + \sum_{n=4}^{\infty} n \left(\frac{44}{48}\right)^2 \left(\frac{48}{52}\right)^{n-4} \frac{4}{52} \\ &= \sum_{n=1}^2 n \left(\frac{44}{48}\right)^{n-1} \frac{4}{48} + \sum_{m=0}^{\infty} (m+4) \left(\frac{44}{48}\right)^2 \left(\frac{48}{52}\right)^m \frac{4}{52}. \end{aligned}$$

But

$$\begin{aligned} \sum_{m=0}^{\infty} (m+4)r^m &= \sum_{m=0}^{\infty} \left(3r^m + \frac{d}{dr} r^{m+1}\right) \\ &= \frac{3}{1-r} + \frac{d}{dr} \left(\frac{r}{1-r}\right) \\ &= \frac{3}{1-r} + \frac{1}{(1-r)^2}, \end{aligned}$$

so

$$\begin{aligned} E[X|Y = 3] &= \frac{4}{48} + 2\left(\frac{44}{48}\right)\left(\frac{4}{48}\right) + \left(\frac{44}{48}\right)^2\left(\frac{4}{52}\right)\left(\frac{3}{1-(48/52)} + \frac{1}{(1-(48/52))^2}\right) \\ &= \frac{4}{48} + 2\left(\frac{44}{48}\right)\left(\frac{4}{48}\right) + \left(\frac{44}{48}\right)^2\left(\frac{4}{52}\right)\left(\frac{3 \times 52}{4} + \frac{52^2}{4^2}\right) \\ &= \frac{1}{12} + 2\left(\frac{11}{12}\right)\left(\frac{1}{12}\right) + 3\left(\frac{11}{12}\right)^2 + \frac{52}{4}\left(\frac{11}{12}\right)^2 \\ &= \frac{985}{72} \approx 13.68. \end{aligned}$$

Given that the first ace appeared on draw 3, it takes an average of between 13 and 14 draws until the first king appears. Compare this to the unconditional $E[X]$. Since $X \sim \text{Geometric}(\frac{4}{52})$ we know $E[X] = \frac{52}{4} = 13$. In other words, on average it takes 13 draws to observe the first king. But given that the first ace appeared on draw three, we should expect to need about 0.68 draws more (on average) to see the first king. ■

Recall how Theorem 1.3.2 described a way in which a non-conditional probability could be calculated in terms of conditional probabilities. There is an analogous theorem for expected value.

THEOREM 4.4.4. Let $X : S \rightarrow T$ be a discrete random variable and let $\{B_i : i \geq 1\}$ be a disjoint collection of events for which $P(B_i) > 0$ for all i and such that $\bigcup_{i=1}^{\infty} B_i = S$. Suppose $P(B_i)$ and $E[X|B_i]$ are known. Then $E[X]$ may be computed as

$$E[X] = \sum_{i=1}^{\infty} E[X|B_i]P(B_i).$$

Proof - Using Theorem 1.3.2 and the definition of conditional expectation,

$$\begin{aligned} \sum_{i=1}^{\infty} E[X|B_i]P(B_i) &= \sum_{i=1}^{\infty} \sum_{t \in T} t \cdot P(X = t|B_i)P(B_i) \\ &= \sum_{t \in T} \sum_{i=1}^{\infty} t \cdot P(X = t|B_i)P(B_i) \\ &= \sum_{t \in T} t \cdot P(X = t) = E[X]. \end{aligned}$$

■

EXAMPLE 4.4.5. A venture capitalist estimates that regardless of whether the economy strengthens, weakens, or remains the same in the next fiscal quarter, a particular investment could either gain or lose money. However, he figures that if the economy strengthens, the investment should, on average, earn 3 million dollars. If the economy remains the same, he figures the expected gain on the investment will be 1 million dollars, while if the economy weakens, the investment will, on average, lose 1 million dollars. He also trusts economic forecasts which predict a 50% chance of a weaker economy, a 40% chance of a stagnant economy, and a 10% chance of a stronger economy. What should he calculate is the expected return on the investment?

Let X be the return on investment and let A , B , and C represent the events that the economy will be stronger, the same, and weaker in the next quarter, respectively. Then the estimates on return give the following information in millions:

$$E[X|A] = 3; \quad E[X|B] = 1; \quad \text{and} \quad E[X|C] = -1.$$

Therefore,

$$\begin{aligned} E[X] &= E[X|A]P(A) + E[X|B]P(B) + E[X|C]P(C) \\ &= 3(0.1) + 1(0.4) + (-1)(0.5) = 0.2 \end{aligned}$$

The expected return on investment is \$200,000. ■

When the conditioning event is described in terms of outcomes of a random variable, Theorem 4.4.4 can be written in another useful way.

THEOREM 4.4.6. Let X and Y be two discrete random variables on a sample space S with $Y : S \rightarrow T$. Let $g : T \rightarrow \mathbb{R}$ be defined as $g(y) = E[X|Y = y]$. Then

$$E[g(Y)] = E[X].$$

It is common to use $E[X|Y]$ to denote $g(Y)$ after which the theorem may be expressed as $E[E[X|Y]] = E[X]$. This can be slightly confusing notation, but one must keep in mind that the exterior expected value in the expression $E[E[X|Y]]$ refers to the average of $E[X|Y]$ viewed as a function of Y .

Proof - As y ranges over T , the events $(Y = y)$ are disjoint and cover all of S . Therefore, by Theorem 4.4.4,

$$\begin{aligned} E[g(Y)] &= \sum_{y \in T} g(y)P(Y = y) \\ &= \sum_{y \in T} E[X|Y = y]P(Y = y) \\ &= E[X]. \end{aligned}$$

■

EXAMPLE 4.4.7. Let $Y \sim \text{Uniform}(\{1, 2, \dots, n\})$ and let X be the number of heads on Y flips of a coin. What is the expected value of X ?

Without Theorem 4.4.6 this problem would require computing many complicated probabilities. However, it is made much simpler by noting that the distribution of X is given conditionally by $(X|Y = j) \sim \text{Binomial}(j, \frac{1}{2})$. Therefore we know $E[X|Y = j] = \frac{j}{2}$. Using the notation above, this may be written as $E[X|Y] = \frac{Y}{2}$ after which

$$E[X] = E[E[X|Y]] = E\left[\frac{Y}{2}\right] = \frac{1}{2} \frac{n+1}{2} = \frac{n+1}{4}.$$

■

Though it requires a somewhat more complicated formula, the variance of a random variable can be computed from conditional information.

THEOREM 4.4.8. Let $X : S \rightarrow T$ be a discrete random variable and let $\{B_i : i \geq 1\}$ be a disjoint collection of events for which $P(B_i) > 0$ for all i and such that $\bigcup_{i=1}^{\infty} B_i = S$. Suppose $E[X|B_i]$ and $\text{Var}[X|B_i]$ are known. Then $\text{Var}[X]$ may be computed as

$$\text{Var}[X] = \left(\sum_{i=1}^{\infty} (\text{Var}[X|B_i] + (E[X|B_i])^2) P(B_i) \right) - (E[X])^2.$$

Proof- First note that $\text{Var}[X|B_i] = E[X^2|B_i] - (E[X|B_i])^2$, and so

$$\text{Var}[X|B_i] + (E[X|B_i])^2 = E[X^2|B_i].$$

Therefore,

$$\sum_{i=1}^{\infty} (\text{Var}[X|B_i] + (E[X|B_i])^2) P(B_i) = \sum_{i=1}^{\infty} E[X^2|B_i] P(B_i),$$

but the right hand side of this equation is $E[X^2]$ from Theorem 4.4.4. The fact that $\text{Var}[X] = E[X^2] - (E[X])^2$ completes the proof of the theorem. ■

As with expected value, this formula may be rewritten in a different form if the conditioning events describe the outcomes of a random variable.

THEOREM 4.4.9. Let X and $Y : S \rightarrow T$ be two discrete random variables on a sample space S . As in Theorem 4.4.6 let $g(y) = E[X|Y = y]$. Let $h(y) = \text{Var}[X|Y = y]$. Denoting $g(Y)$ by $E[X|Y]$ and denoting $h(Y)$ by $\text{Var}[X|Y]$, then

$$\text{Var}[X] = E[\text{Var}[X|Y]] + \text{Var}[E[X|Y]].$$

Proof - First consider the following three facts:

- (1) $\sum_{t \in T} \text{Var}[X|Y = t]P(Y = t) = E[\text{Var}[X|Y]]$;
- (2) $\sum_{t \in T} (E[X|Y = t])^2 P(Y = t) = E[(E[X|Y])^2]$; and
- (3) $\text{Var}[E[X|Y]] = E[(E[X|Y])^2] - (E[E[X|Y]])^2 = E[(E[X|Y])^2] - (E[X])^2$.

Then from Theorem 4.4.8,

$$\begin{aligned} \text{Var}[X] &= \sum_{t \in T} (\text{Var}[X|Y = t] + (E[X|Y = t])^2)P(Y = t) - (E[X])^2 \\ &= \sum_{t \in T} \text{Var}[X|Y = t]P(Y = t) + \sum_{t \in T} (E[X|Y = t])^2 P(Y = t) - (E[X])^2 \\ &= E[\text{Var}[X|Y]] + E[(E[X|Y])^2] - (E[X])^2 \\ &= E[\text{Var}[X|Y]] + \text{Var}[E[X|Y]]. \end{aligned}$$

■

EXAMPLE 4.4.10. The number of eggs N found in nests of a certain species of turtles has a Poisson distribution with mean λ . Each egg has probability p of being viable and this event is independent from egg to egg. Find the mean and variance of the number of viable eggs per nest.

Let N be the total number of eggs in a nest and X the number of viable ones. Then if $N = n$, X has a binomial distribution with number of trials n and probability p of success for each trial. Thus, if $N = n$, X has mean np and variance $np(1 - p)$. That is,

$$E[X|N = n] = np; \quad \text{Var}[X|N = n] = np(1 - p)$$

or

$$E[X|N] = pN; \quad \text{Var}[X|N] = p(1 - p)N.$$

Hence

$$E[X] = E[E[X|N]] = E[pN] = pE[N] = p\lambda$$

and

$$\begin{aligned} \text{Var}[X] &= E[\text{Var}[X|N]] + \text{Var}[E[X|N]] \\ &= E[p(1 - p)N] + \text{Var}[pN] = p(1 - p)E[N] + p^2\text{Var}[N]. \end{aligned}$$

Since N is Poisson we know that $E[N] = \text{Var}[N] = \lambda$, so that

$$E[X] = p\lambda \quad \text{and} \quad \text{Var}[X] = p(1 - p)\lambda + p^2\lambda = p\lambda.$$

■

EXERCISES

Ex. 4.4.1. Let $X \sim \text{Geometric}(p)$ and let A be event $(X \leq 3)$. Calculate $E[X|A]$ and $\text{Var}[X|A]$.

Ex. 4.4.2. Calculate the variance of the quantity X from Example 4.4.7.

Ex. 4.4.3. Return to Example 4.4.5. Suppose that, in addition to the estimates on average return, the investor had estimates on the standard deviations. If the economy strengthens or weakens, the estimated standard deviation is 3 million dollars, but if the economy stays the same, the estimated standard deviation is 2 million dollars. So, in millions of dollars,

$$SD[X|A] = 3; \quad SD[X|B] = 2; \quad \text{and} \quad SD[X|C] = 3.$$

Use this information, together with the conditional expectations from Example 4.4.5 to calculate $Var[X]$.

Ex. 4.4.4. A standard light bulb has an average lifetime of four years with a standard deviation of one year. A Super D-Lux lightbulb has an average lifetime of eight years with a standard deviation of three years. A box contains many bulbs – 90% of which are standard bulbs and 10% of which are Super D-Lux bulbs. A bulb is selected at random from the box. What are the average and standard deviation of the lifetime of the selected bulb?

Ex. 4.4.5. Let X and Y be described by the joint distribution

	$X = -1$	$X = 0$	$X = 1$
$Y = -1$	1/15	2/15	2/15
$Y = 0$	2/15	1/15	2/15
$Y = 1$	2/15	2/15	1/15

and answer the following questions.

- (a) Calculate $E[X|Y = -1]$.
- (b) Calculate $Var[X|Y = -1]$.
- (c) Describe the distribution of $E[X|Y]$.
- (d) Describe the distribution of $Var[X|Y]$.

Ex. 4.4.6. Let X and Y be discrete random variables. Let x be in the range of X and let y be in the range of Y .

- (a) Suppose X and Y are independent. Show that $E[X|Y = y] = E[X]$ (and so $E[X|Y] = E[X]$).
- (b) Show that $E[X|X = x] = x$ (and so $E[X|X] = X$). (From results in this section we know $E[X|Y]$ is always a random variable with expected value equal to $E[X]$. The results above in some sense show two extremes. When X and Y are independent, $E[X|Y]$ is a constant random variable $E[X]$. When X and Y are equal, $E[X|X]$ is just X itself).

Ex. 4.4.7. Let $X \sim \text{Uniform } \{1, 2, \dots, n\}$ be independent of $Y \sim \text{Uniform } \{1, 2, \dots, n\}$. Let $Z = \max(X, Y)$ and $W = \min(X, Y)$.

- (a) Find the joint distribution of (Z, W) .
- (b) Find $E[Z | W]$.

4.5 COVARIANCE AND CORRELATION

When faced with two different random variables, we are frequently interested in how the two different quantities relate to each other. Often the purpose of this is to predict something about one variable knowing information about the other. For instance, if rainfall amounts in July affect the quantity of corn harvested in August, then a farmer, or anyone else keenly interested in the supply and demand of the agriculture industry, would like to be able to use the July information to help make predictions about August costs.

4.5.1 Covariance

Just as we developed the concepts of expected value and standard deviation to summarize a single random variable, we would like to develop a number that describes something about how two different random variables X and Y relate to each other.

DEFINITION 4.5.1. (Covariance of X and Y) Let X and Y be two discrete random variables on a sample space S . Then the “covariance of X and Y ” is defined as

$$\text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])]. \quad (4.5.1)$$

Since it is defined in terms of an expected value, there is the possibility that the covariance may be infinite or not defined at all because the sum describing the expectation is divergent.

Notice that if X is larger than its average at the same time that Y is larger than its average (or if X is smaller than its average at the same time Y is smaller than its average) then $(X - E[X])(Y - E[Y])$ will contribute a positive result to the expected value describing the covariance. Conversely, if X is smaller than $E[X]$ while Y is larger than $E[Y]$ or vice versa, a negative result will be contributed toward the covariance. This means that when two variables tend to be both above average or both below average simultaneously, the covariance will typically be positive (and the variables are said to be positively correlated), but when one variable tends to be above average when the other is below average, the covariance will typically be negative (and the variables are said to be negatively correlated). When $\text{Cov}[X, Y] = 0$ the variables X and Y are said to be “uncorrelated”.

For example, suppose X and Y are the height and weight, respectively, of an individual randomly selected from a large population. We might expect that $\text{Cov}[X, Y] > 0$ since people who are taller than average also tend to be heavier than average and people who are shorter than average tend to be lighter. Conversely suppose X and Y represent elevation and air density at a randomly selected point on Earth. We might expect $\text{Cov}[X, Y] < 0$ since locations at a higher elevation tend to have thinner air.

EXAMPLE 4.5.2. Consider a pair of random variables X and Y with joint distribution

	$X = -1$	$X = 0$	$X = 1$
$Y = -1$	1/15	2/15	2/15
$Y = 0$	2/15	1/15	2/15
$Y = 1$	2/15	2/15	1/15

By a routine calculation of the marginal distributions it can be shown that $X, Y \sim \text{Uniform}(\{-1, 0, 1\})$ and therefore that $E[X] = E[Y] = 0$. However, it is clear from the joint distribution that when

$X = -1$, then Y is more likely to be above average than below, while when $X = 1$, then Y is more likely to be below average than above. This suggests the two random variables should have a negative correlation. In fact, we can calculate

$$E[XY] = (-1)\left(\frac{4}{15}\right) + 0\left(\frac{9}{15}\right) + 1\left(\frac{2}{15}\right) = -\frac{2}{15},$$

and therefore $Cov[X, Y] = E[XY] - E[X]E[Y] = -\frac{2}{15}$. ■

As its name suggests, the covariance is closely related to the variance.

THEOREM 4.5.3. *Let X be a discrete random variable. Then*

$$Cov[X, X] = Var[X].$$

Proof - $Cov[X, X] = E[(X - E[X])(X - E[X])] = E[(X - E[X])^2] = Var[X]$. ■

With Theorem 4.2.5 it was shown that $Var[X] = E[X^2] - (E[X])^2$, which provided an alternate formula for the variance. There is an analogous alternate formula for the covariance.

THEOREM 4.5.4. *Let X and Y be discrete random variables with finite mean for which $E[XY]$ is also finite. Then*

$$Cov[X, Y] = E[XY] - E[X]E[Y].$$

Proof - Using the linearity properties of expected value,

$$\begin{aligned} Cov[X, Y] &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY - XE[Y] - E[X]Y + E[X]E[Y]] \\ &= E[XY] - E[XE[Y]] - E[E[X]Y] + E[E[X]E[Y]] \\ &= E[XY] - E[Y]E[X] - E[X]E[Y] + E[X]E[Y] \\ &= E[XY] - E[X]E[Y]. \end{aligned}$$
■

As with the expected value, the covariance is a linear quantity. It is also related to the concept of independence.

THEOREM 4.5.5. *Let X , Y , and Z be discrete random variables, and let $a, b \in \mathbb{R}$. Then,*

- (a) $Cov[X, Y] = Cov[Y, X]$;
- (b) $Cov[X, aY + bZ] = a \cdot Cov[X, Y] + b \cdot Cov[X, Z]$;
- (c) $Cov[aX + bY, Z] = a \cdot Cov[X, Z] + b \cdot Cov[Y, Z]$; and
- (d) If X and Y are independent with a finite covariance, then $Cov[X, Y] = 0$.

Proof of (1) - This follows immediately from the definition.

$$\begin{aligned} Cov[X, Y] &= E[(X - E[X])(Y - E[Y])] \\ &= E[(Y - E[Y])(X - E[X])] = Cov[Y, X]. \end{aligned}$$

Therefore, reversing the roles of X and Y does not change the correlation.

Proof of (2) - This follows from linearity properties of expected value. Using Theorem 4.5.4

$$\begin{aligned} \text{Cov}[X, aY + bZ] &= E[X(aY + bZ)] - E[X]E[aY + bZ] \\ &= a \cdot E[XY] + b \cdot E[XZ] - a \cdot E[X]E[Y] - b \cdot E[X]E[Z] \\ &= a \cdot (E[XY] - E[X]E[Y]) + b \cdot (E[XZ] - E[X]E[Z]) \\ &= a \cdot \text{Cov}[X, Y] + b \cdot \text{Cov}[X, Z] \end{aligned}$$

Proof of (3) - This proof is essentially the same as that of (2) and is left as an exercise.

Proof of (4) - We have previously seen that if X and Y are independent, then $E[XY] = E[X]E[Y]$. Using Theorem 4.5.4 it follows that

$$\text{Cov}[X, Y] = E[XY] - E[X]E[Y] = 0.$$

Though independence of X and Y guarantees that they are uncorrelated, the converse is not true. It is possible that $\text{Cov}[X, Y] = 0$ and yet that X and Y are dependent, as the next example shows.

EXAMPLE 4.5.6. Let X, Y be two discrete random variables taking values $\{-1, 1\}$. Suppose their joint distribution $P(X = x, Y = y)$ is given by the table

	x=-1	x=1
y=-1	0.3	0.2
y=1	0.3	0.2

By summing the columns and rows respectively,

$$P(X = 1) = 0.4 \text{ and } P(X = -1) = 0.6, \text{ while}$$

$$P(Y = 1) = 0.5 \text{ and } P(Y = -1) = 0.5.$$

Moreover,

$$\begin{aligned} E[XY] &= (1)(-1)P(X = 1, Y = -1) + (-1)(1)P(X = -1, Y = 1) \\ &\quad + (1)(1)P(X = 1, Y = 1) + (-1)(-1)P(X = -1, Y = -1) \\ &= -0.3 - 0.2 + 0.2 + 0.3 = 0, \\ E[X] &= (1)0.4 + (-1)0.6 = -0.2, \\ E[Y] &= (1)0.5 + (-1)0.5 = 0, \end{aligned}$$

implying that $\text{Cov}[X, Y] = 0$. As

$$P(X = 1, Y = 1) = 0.2 \neq 0.1 = P(X = 1)P(Y = 1),$$

they are not independent random variables. ■

4.5.2 Correlation

The possible size of $\text{Cov}[X, Y]$ has upper and lower bounds based on the standard deviations of the two variables.

THEOREM 4.5.7. Let X and Y be two discrete random variables both with finite variance. Then

$$-\sigma_X\sigma_Y \leq \text{Cov}[X, Y] \leq \sigma_X\sigma_Y,$$

and therefore $-1 \leq \frac{\text{Cov}[X, Y]}{\sigma_X\sigma_Y} \leq 1$.

Proof - Standardize both variables and consider the expected value of their sum squared. Since this is the expected value of a non-negative quantity,

$$\begin{aligned} 0 &\leq E\left[\left(\frac{X - \mu_X}{\sigma_X} + \frac{Y - \mu_Y}{\sigma_Y}\right)^2\right] \\ &= E\left[\frac{(X - \mu_X)^2}{\sigma_X^2} + 2\frac{(X - \mu_X)(Y - \mu_Y)}{\sigma_X\sigma_Y} + \frac{(Y - \mu_Y)^2}{\sigma_Y^2}\right] \\ &= \frac{E[(X - \mu_X)^2]}{\sigma_X^2} + \frac{2E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X\sigma_Y} + \frac{E[(Y - \mu_Y)^2]}{\sigma_Y^2} \\ &= 1 + 2\frac{\text{Cov}[X, Y]}{\sigma_X\sigma_Y} + 1. \end{aligned}$$

Solving the inequality for the covariance yields

$$\text{Cov}[X, Y] \geq -\sigma_X\sigma_Y.$$

A similar computation (see Exercises) for the expected value of the squared difference of the standardized variables shows

$$\text{Cov}[X, Y] \leq \sigma_X\sigma_Y.$$

Putting both inequalities together proves the theorem. ■

DEFINITION 4.5.8. The quantity $\frac{\text{Cov}[X, Y]}{\sigma_X\sigma_Y}$ from Theorem 4.5.7 is known as the “correlation” of X and Y and is often denoted as $\rho[X, Y]$. Thinking in terms of dimensional analysis, both the numerator and denominator include the units of X and the units of Y . The correlation, therefore, has no units associated with it. It is thus a dimensionless rescaling of the covariance and is frequently used as an absolute measure of trends between the two variables.

EXERCISES

Ex. 4.5.1. Consider the experiment of flipping two coins. Let X be the number of heads among the coins and let Y be the number of tails among the coins.

- (a) Should you expect X and Y to be positively correlated, negatively correlated, or uncorrelated? Why?
- (b) Calculate $\text{Cov}[X, Y]$ to confirm your answer to (a).

Ex. 4.5.2. Let $X \sim \text{Uniform}(\{0, 1, 2\})$ and let Y be the number of heads in X flips of a coin.

- (a) Should you expect X and Y to be positively correlated, negatively correlated, or uncorrelated? Why?

(b) Calculate $Cov[X, Y]$ to confirm your answer to (a).

Ex. 4.5.3. Prove part (3) of Theorem 4.5.5.

Ex. 4.5.4. Prove the missing inequality from the proof of Theorem 4.5.7. Specifically, use the inequality

$$0 \leq E\left[\left(\frac{X - \mu_X}{\sigma_X} - \frac{Y - \mu_Y}{\sigma_Y}\right)^2\right]$$

to prove that $Cov[X, Y] \leq \sigma_X \sigma_Y$.

Ex. 4.5.5. Prove that the inequality of Theorem 4.5.7 is an equality if and only if there are $a, b \in \mathbb{R}$ with $a \neq 0$ for which $P(Y = aX + b) = 1$. (Put another way, the correlation of X and Y is ± 1 exactly when Y can be expressed as a non-trivial linear function of X).

Ex. 4.5.6. In previous sections it was shown that if X and Y are independent, then $Var[X + Y] = Var[X] + Var[Y]$. If X and Y are dependent, the result is typically not true, but the covariance provides a way relate the variances of X and Y to the variance of their sum.

(a) Show that for any discrete random variables X and Y ,

$$Var[X + Y] = Var[X] + Var[Y] + 2Cov[X, Y].$$

(b) Use (a) to conclude that when X and Y are positively correlated, then $Var[X + Y] > Var[X] + Var[Y]$, while when X and Y are negatively correlated, $Var[X + Y] < Var[X] + Var[Y]$.

(c) Suppose X_i $1 \leq i \leq n$ are discrete random variables with finite variance and covariances. Use induction and (a) to conclude that

$$Var\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n Var[X_i] + 2 \sum_{1 \leq i < j \leq n} Cov[X_i, X_j].$$

4.6 EXCHANGEABLE RANDOM VARIABLES

We conclude this section with a discussion on exchangeable random variables. In brief we say that a collection of random variables is exchangeable if the joint probability mass function of (X_1, X_2, \dots, X_n) is a symmetric function. In other words, the distribution of (X_1, X_2, \dots, X_n) is independent of the order in which the X_i 's appear. In particular any collection of mutually independent random variables is exchangeable.

DEFINITION 4.6.1. Let $n \geq 2$ and $\sigma : \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, n\}$ be a bijection. We say that a subset T of \mathbb{R}^n is symmetric if

$$(x_1, x_2, \dots, x_n) \in T \iff (x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(n)}) \in T$$

for all $(x_1, x_2, \dots, x_n) \in \mathbb{R}^n$. For any symmetric set T , a function $f : T \rightarrow \mathbb{R}$ is symmetric if

$$f(x_1, x_2, \dots, x_n) = f(x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(n)})$$

for all $(x_1, x_2, \dots, x_n) \in \mathbb{R}^n$.

A bijection $\sigma : \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, n\}$ is often referred to as a permutation of $\{1, 2, \dots, n\}$. When $n = 2$ the function f would be symmetric if $f(x, y) = f(y, x)$ for all $x, y \in \mathbb{R}$.

DEFINITION 4.6.2. Let $n \geq 1$ and X_1, X_2, \dots, X_n be discrete random variables. We say that X_1, X_2, \dots, X_n is a collection of exchangeable random variables if the joint probability mass function given by

$$f(x_1, x_2, \dots, x_n) = P(X_1 = x_1, \dots, X_n = x_n)$$

is a symmetric function.

In particular, X_1, X_2, \dots, X_n are exchangeable then for any one of the possible $n!$ permutations, σ , of $\{1, 2, \dots, n\}$, X_1, X_2, \dots, X_n and $X_{\sigma(1)}, X_{\sigma(2)}, \dots, X_{\sigma(n)}$ have the same distribution.

EXAMPLE 4.6.3. Suppose we have an urn of m distinct objects labelled $\{1, 2, \dots, m\}$. Objects are drawn at random from the urn without replacements till the urn is empty. Let X_i be the label of the i -th object that is drawn. Then X_1, X_2, \dots, X_m is a particular ordering of the objects in the urn. Since each ordering is equally likely and there are $m!$ possible orderings we have that the joint probability mass function

$$f(x_1, x_2, \dots, x_m) = P(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m) = \frac{1}{m!},$$

whenever $x_i \in \{1, 2, \dots, m\}$ with $x_i \neq x_j$. As the function is a constant function on the symmetric set $\{1, 2, \dots, m\}$, it is clearly symmetric. So the random variables X_1, X_2, \dots, X_m are exchangeable. ■

THEOREM 4.6.4. Let X_1, X_2, \dots, X_n be a collection of exchangeable random variables on a sample space S . Then for any $i, j \in \{1, 2, \dots, n\}$, X_i and X_j have the same marginal distribution.

Proof - The random variables (X_1, X_2, \dots, X_n) are exchangeable. Then we have for any permutation σ and $x_i \in \text{Range}(X_i)$

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = P(X_{\sigma(1)} = x_1, X_{\sigma(2)} = x_2, \dots, X_{\sigma(n)} = x_n).$$

As this is true for all permutations σ all the random variables must have same range. Otherwise if any two of them differ the we could get a contradiction by choosing an appropriate permutation.

Let T denote the common range. Let $i \in \{2, \dots, n\}, a, b \in T$. Let

$$A = \{x_j \in T : 1 \leq j \neq i, 1 \leq j \leq n\}$$

By using the exchangeable property with the permutation σ that is given by $\sigma(i) = 1, \sigma(1) = i$ and $\sigma(j) = j$ for all $j \neq 1, i$. We have that for any $x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n \in A$

$$\begin{aligned} & P(X_1 = a, X_2 = x_2, \dots, X_{i-1} = x_{i-1}, X_i = b, X_{i+1} = x_{i+1}, \dots, X_n = x_n) \\ &= P(X_1 = b, X_2 = x_2, \dots, X_{i-1} = x_{i-1}, X_i = a, X_{i+1} = x_{i+1}, \dots, X_n = x_n). \end{aligned}$$

Therefore,

$$\begin{aligned}
P(X_1 = a) &= P\left(\bigcup_{b \in T} X_1 = a, X_i = b\right) \\
&= \sum_{b \in T} P(X_1 = a, X_i = b) \\
&= \sum_{b \in T} P\left(\bigcup_{x_j \in A} X_1 = a, X_2 = x_2, \dots, X_{i-1} = x_{i-1}, X_i = b, X_{i+1} = x_{i+1}, \dots, X_n = x_n\right) \\
&= \sum_{b \in T} \sum_{x_j \in A} P(X_1 = a, X_2 = x_2, \dots, X_i = b, \dots, X_n = x_n) \\
&= \sum_{b \in T} \sum_{x_j \in A} P(X_1 = b, X_2 = x_2, \dots, X_i = a, \dots, X_n = x_n) \\
&= \sum_{b \in T} \sum_{x_j \in A} P\left(\bigcup_{x_j \in A} X_1 = b, X_2 = x_2, \dots, X_{i-1} = x_{i-1}, X_i = a, X_{i+1} = x_{i+1}, \dots, X_n = x_n\right) \\
&= \sum_{b \in T} P(X_1 = b, X_i = a) \\
&= P\left(\bigcup_{b \in T} X_1 = b, X_i = a\right) \\
&= P(X_i = a)
\end{aligned}$$

So the distribution of X_i is the same as the distribution of X_1 and hence all of them have the same distribution. ■

EXAMPLE 4.6.5. (Sampling without Replacement) An urn contains b black balls and r red balls. A ball is drawn at random and its colour noted. This procedure is repeated n times. Assume that $n \leq b + r$. Let $\max 0, n - r \leq k \leq \min(n, b)$. In this example we examine the random variables X_i given by

$$X_i = \begin{cases} 1 & \text{if } i\text{-th ball drawn is black} \\ 0 & \text{otherwise} \end{cases}$$

We have already seen that (See Theorem 2.3.2 and Example 2.3.1)

$$P(k \text{ black balls are drawn in } n \text{ draws}) = \binom{n}{k} \frac{\prod_{i=0}^{k-1} (b-i) \prod_{i=0}^{m-k-1} (r-i)}{\prod_{i=0}^{m-1} (r+b-i)}.$$

Using the same proof we see that the joint probability mass function of (X_1, X_2, \dots, X_n) is given by

$$f(x_1, x_2, \dots, x_n) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = \frac{\prod_{i=0}^n x_i^{-1} (b-i) \prod_{i=0}^n x_i^{-k-1} (r-i)}{\prod_{i=0}^n (r+b-i)},$$

where $x_i \in \{0, 1\}$. It is clear from the right hand side of the above that the function f depends only on the $\sum_{i=1}^n x_i$. Hence any permutation of the x_i 's will not change the value of f . So f is a symmetric function and the random variables are exchangeable. Therefore, by Theorem 4.6.4 we know that for any $1 \leq i \leq n$,

$$P(X_i = 1) = P(X_1 = 1) = \frac{b}{b+r}.$$

So we can conclude that they are all identically distributed as Bernoulli ($\frac{b}{b+r}$) and the probability of choosing a black ball in the i -th draw is $\frac{b}{b+r}$ (See Exercise 4.6.4 for a similar result). Further for any i, j

$$\begin{aligned} \text{Cov}[X_i, X_j] &= E[X_i X_j] - E[X_i]E[X_j] \\ &= E[X_1 X_2] - \left(\frac{b}{b+r}\right)^2 \\ &= \frac{b(b-1)}{(b+r)(b+r-1)} - \left(\frac{b}{b+r}\right)^2 \\ &= \frac{-br}{(b+r)^2(b+r-1)} \end{aligned}$$

Finally, we observe that $Y = \sum_{i=1}^n X_i$ is a Hypergeometric $(b+r, b, m)$. Exchangeability thus provides another alternative way to compute the mean and variance of Y . Using the linearity of expectation provided by Theorem 4.1.7, we have

$$E[Y] = E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i] = n \frac{b}{b+r}.$$

and by Exercise 4.5.6,

$$\begin{aligned} \text{Var}[Y] &= \text{Var}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \text{Var}[X_i] + 2 \sum_{1 \leq i < j \leq n} \text{Cov}[X_i, X_j] \\ &= n \text{Var}[X_1] + n(n-1) \text{Cov}[X_1, X_2] \\ &= n \frac{br}{(b+r)^2} + n(n-1) \left(\frac{-br}{(b+r)^2(b+r-1)}\right) \\ &= n \frac{br}{(b+r)^2} \frac{b+r-n}{b+r-1}. \end{aligned}$$



EXERCISES

Ex. 4.6.1. Suppose X_1, X_2, \dots, X_n are exchangeable random variables. For any $2 \leq m < n$, show that X_1, X_2, \dots, X_m are also a collection of exchangeable random variables.

Ex. 4.6.2. Suppose X_1, X_2, \dots, X_n are exchangeable random variables. Let T denote their common range. Suppose $b : T \rightarrow \mathbb{R}$. Show that $b(X_1), b(X_2), \dots, b(X_n)$ is also a collection of exchangeable random variables.

Ex. 4.6.3. Suppose n cards are drawn from a standard pack of 52 cards without replacement (so we will assume $n \leq 52$). For $1 \leq i \leq n$, let X_i be random variables given by

$$X_i = \begin{cases} 1 & \text{if } i\text{-th card drawn is black in colour} \\ 0 & \text{otherwise} \end{cases}$$

- (a) Suppose $n = 52$. Using Example 4.6.3 and the Exercise 4.6.2 show that $(X_1, X_2, X_3, \dots, X_n)$ are exchangeable.
- (b) Show that $(X_1, X_2, X_3, \dots, X_n)$ are exchangeable for any $2 \leq n \leq 52$. Hint: If $n < 52$ extend the sample to exhaust the deck of cards. Use (a) and Exercise 4.6.1

- (c) Find the probability that the second and fourth card drawn have the same colour.

Ex. 4.6.4. (Polya Urn Scheme) An urn contains b black balls and r red balls. A ball is drawn at random and its colour noted. Then it is replaced along with $c \geq 0$ balls of the same colour. This procedure is repeated n times.

- (a) Let $1 \leq k \leq m \leq n$. Show that

$$P(k \text{ black balls are drawn in } m \text{ draws}) = \binom{m}{k} \frac{\prod_{i=0}^{k-1} (b + ci) \prod_{i=0}^{m-k-1} (r + ci)}{\prod_{i=0}^{m-1} (r + b + ci)}$$

- (b) Let $1 \leq i \leq n$ and random variables X_i be given by

$$X_i = \begin{cases} 1 & \text{if } i\text{-th ball drawn is black} \\ 0 & \text{otherwise} \end{cases}$$

Show that the collection of random variables is exchangeable.

- (c) Let $1 \leq m \leq n$. Let B_m be the event that the m -th ball drawn is black. Show that

$$P(B_m) = \frac{b}{b+r}.$$

