

# Modeling with Probability

Peter Howard

Spring 2022

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Definitions and Axioms</b>	<b>3</b>
<b>3</b>	<b>Counting Arguments</b>	<b>6</b>
3.1	Permutations . . . . .	6
3.2	Combinations . . . . .	7
3.3	Combinations with Repetition . . . . .	11
3.4	Multiple Combinations . . . . .	12
<b>4</b>	<b>Conditional Probability</b>	<b>14</b>
4.1	Computing Probabilities by Conditioning . . . . .	16
4.2	Bayes' Lemma . . . . .	19
4.3	Independent Events . . . . .	21
<b>5</b>	<b>Discrete Random Variables</b>	<b>22</b>
5.1	Expected Value . . . . .	23
5.1.1	Properties of Expected Value . . . . .	26
5.1.2	Conditional Expected Value . . . . .	28
5.2	Variance and Covariance . . . . .	30
5.3	Cumulative Distribution Functions . . . . .	31
5.4	Probability Mass Functions . . . . .	32
<b>6</b>	<b>Game Theory</b>	<b>34</b>
6.1	Zero-Sum Games . . . . .	34
6.1.1	Dominant Strategies . . . . .	35
6.1.2	Saddle Points . . . . .	38
6.1.3	Minimax Method . . . . .	40
6.1.4	Mixed Strategies . . . . .	41
6.1.5	The Method of Oddments . . . . .	44
6.1.6	The Minimax Theorem . . . . .	45

<b>7</b>	<b>Continuous Random Variables</b>	<b>51</b>
7.1	Cumulative Distribution Functions . . . . .	51
7.2	Probability Density Functions . . . . .	52
7.3	Expected Value, Variance, and Covariance . . . . .	53
7.4	Identifying Probability Density Functions . . . . .	54
7.5	Useful Probability Density Functions . . . . .	55
7.6	More Probability Density Functions . . . . .	65
7.7	Joint Probability Density Functions . . . . .	68
<b>8</b>	<b>Maximum Likelihood Estimators</b>	<b>69</b>
8.1	Maximum Likelihood Estimation for Discrete Random Variables . . . . .	69
8.2	Maximum Likelihood Estimation for Continuous Random Variables . . . . .	70
<b>9</b>	<b>Simulating a Random Process</b>	<b>73</b>
9.1	Simulating Uniform Random Variables . . . . .	75
9.2	Simulating Discrete Random Variables . . . . .	75
9.3	Simulating Gaussian Random Variables . . . . .	76
9.4	Simulating More General Random Variables . . . . .	76
9.4.1	The Rejection Method. . . . .	77
9.5	Application to Queuing Theory . . . . .	80
9.6	Limit Theorems . . . . .	81
<b>10</b>	<b>Hypothesis Testing</b>	<b>92</b>
10.1	General Hypothesis Testing . . . . .	92
10.2	Hypothesis Testing for Distributions . . . . .	94
10.2.1	Empirical Distribution Functions . . . . .	94
<b>11</b>	<b>Application to Finance</b>	<b>99</b>
11.1	Random Walks . . . . .	99
11.2	Brownian Motion . . . . .	100
11.3	Stochastic Differential Equations . . . . .	103

“We see that the theory of probability is at bottom only common sense reduced to calculation; it makes us appreciate with exactitude what reasonable minds feel by a sort of instinct, often without being able to account for it.... It is remarkable that a science which began with the consideration of games of chance should have become the most important object of human knowledge.... The most important questions of life are, for the most part, really only problems of probability.”

Pierre Simon, Marquis de Laplace, *Théorie Analytique des Probabilités*, 1812

“The gambling passion lurks at the bottom of every heart.”

Honoré de Balzac

# 1 Introduction

Although games of chance have been around in one form or another for thousands of years, the first person to attempt the development of a systematic theory for such games seems to have been the Italian mathematician, physician, astrologer and—yes—gambler Gerolamo Cardano (1501–1576). Cardano is perhaps best known for his study of cubic and quartic algebraic equations, which he solved in his 1545 text *Ars Magna*—solutions which required his keeping track of  $\sqrt{-1}$ . He did not develop a theory of complex numbers, but is largely regarded as the first person to recognize the possibility of using what has now become the theory of complex numbers. He is also remembered as an astrologer who made many bold predictions, including a horoscope of Jesus Christ (1554), and for having predicted the precise day on which he would die and then (or at least as the story goes) committing suicide on that day.<sup>1</sup>

In 1654, Antoine Gombaud Chevalier de Mere, a French nobleman and professional gambler called Blaise Pascal’s (1623–1662) attention to a curious game of chance: was it worthwhile betting even money (original bet is either doubled or lost) that double sixes would turn up at least once in 24 throws of a pair of fair dice. This led to a long correspondence between Pascal and Pierre de Fermat (1601–1665, of Fermat’s Last Theorem fame) in which the fundamental principles of probability theory were formulated for the first time. In these notes we will review some of their key observations, along with elements of more modern theory.

## 2 Definitions and Axioms

Suppose we flip a fair coin twice and record each time whether it lands heads (H) or tails (T). The list of all possible outcomes for this experiment is

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

which constitutes a set that we refer to as the *sample space* for this experiment. Each member of  $S$  is referred to as an *outcome*. In this case, finding the probability of any particular outcome is straightforward. For example, the probability of getting two heads, which we will denote  $P(HH)$ , is  $P(HH) = 1/4$ . Any subset,  $E$ , of the sample space is an *event*. In the example above,  $E = \{(HH), (HT)\}$  is the event that heads appears on the first flip, and  $F = \{(HT), (TT)\}$  is the event that tails appears on the second flip.

**Definition 2.1.** *For any two sets (events)  $A$  and  $B$ , we define the following:*

1. (*Intersection*)  $A \cap B =$  all outcomes in both  $A$  and  $B$  (in our example,  $E \cap F = \{(HT)\}$ ).
2. (*Union*)  $A \cup B =$  all outcomes in either  $A$  or  $B$  (or both) (in our example,  $E \cup F = \{(HH), (HT), (TT)\}$ ).
3. (*Set subtraction*)  $A \setminus B =$  all outcomes in  $A$  that are not also in  $B$  (in our example,  $E \setminus F = \{(HH)\}$ ).

---

<sup>1</sup>Now that’s dedication to your profession.

4. (Complement)  $A^c =$  all outcomes in  $S$  but not in  $A$  (in our example, the complement of  $E$  is  $E^c = \{(TH), (TT)\}$ ).
5. (Containment) If  $x$  is an outcome in  $A$ , we write  $x \in A$ . If every outcome in  $A$  is also an outcome in  $B$ , we write  $A \subset B$ .

**Lemma 2.1.** For any three sets (events)  $A$ ,  $B$ , and  $C$ , the following relations hold:

1.  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$
2.  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$
3.  $(A \cup B)^c = A^c \cap B^c$
4.  $(A \cap B)^c = A^c \cup B^c$

**Proof.** Each of the four items in Lemma 2.1 is proven in a similar way, so let's just consider the first. The strategy is to show containment in two directions,

$$A \cap (B \cup C) \subset (A \cap B) \cup (A \cap C) \quad \text{and} \quad (A \cap B) \cup (A \cap C) \subset A \cap (B \cup C).$$

For the first, we take any  $x \in A \cap (B \cup C)$  and our goal is to show that we must also have  $x \in (A \cap B) \cup (A \cap C)$ . To see this, we observe that the statement  $x \in A \cap (B \cup C)$  asserts that  $x \in A$  and also that  $x$  is in either  $B$  or  $C$ . If  $x \in B$ , then since we already know  $x \in A$  we in fact know that  $x \in A \cap B$ . On the other hand, if  $x \in C$ , then since we already know  $x \in A$  we in fact know that  $x \in A \cap C$ . In this way, we see that  $x$  is either in  $A \cap B$  or in  $A \cap C$ , and this is precisely the statement that  $x \in (A \cap B) \cup (A \cap C)$ .

Turning to the opposite containment, we suppose  $x \in (A \cap B) \cup (A \cap C)$ . I.e., either  $x \in A \cap B$  or  $x \in A \cap C$  (or both), which tells us two things: first, we certainly have  $x \in A$ , and second,  $x$  must be in either  $B$  or  $C$ , meaning that  $x \in B \cup C$ . Since we have both  $x \in A$  and  $x \in B \cup C$ , we can conclude that  $x \in A \cap (B \cup C)$ , which is what we needed to show.  $\square$

One of the first men to systematically develop the theory of probability was Pierre Simon Laplace (1749–1827), who famously said, “At the bottom, the theory of probability is only common sense expressed in numbers.” According to this common sense, we might think of computing the probability that an event  $A$  occurs in something like the following way: if  $n$  experiments are carried out with a sample space  $S$  for which  $A \subset S$  (i.e.,  $A$  is a possible event of the experiment), then we expect that  $P(A)$  is approximately the number of times the event  $A$  occurs divided by  $n$ ,

$$P(A) \cong \frac{\# \text{ of times } A \text{ occurs in the } n \text{ experiments}}{n},$$

and more precisely we might expect this to get better and better as  $n$  increases, so that it's reasonable to write

$$P(A) = \lim_{n \rightarrow \infty} \frac{\# \text{ of times } A \text{ occurs in the } n \text{ experiments}}{n}.$$

In order to develop this intuition into a rigorous mathematical theory, we proceed by assuming a set of axioms that cannot be proven from earlier principles, but which we regard as somehow self-evident.<sup>2</sup>

**Axioms of Probability.** For any sample space  $S$ , we have

**Axiom 1.**  $0 \leq P(E) \leq 1$ , for all events  $E$  in  $S$ .

**Axiom 2.**  $P(S) = 1$ .

**Axiom 3.** If  $E_1, E_2, \dots$  are mutually exclusive events in  $S$  (that is,  $E_k \cap E_j = \emptyset$  for  $k \neq j$ ) then

$$P(E_1 \cup E_2 \cup \dots) = P(E_1) + P(E_2) + \dots$$

We observe that in our simple calculation  $P(HH) = 1/4$ , we have used Axioms 2 and 3 in the following ways. First, we used Axiom 3 to obtain the relation,

$$P((HH) \cup (HT) \cup (TH) \cup (TT)) = P(HH) + P(HT) + P(TH) + P(TT),$$

and second we used Axiom 2 to see that

$$P(HH) + P(HT) + P(TH) + P(TT) = 1,$$

and we finally conclude our calculation by further assuming that each outcome is equally likely.

By combining these axioms with our set relations, we can begin to build up our slate of probability relations. Two examples are given in the next lemma.

**Lemma 2.2.** Suppose  $A$  and  $B$  are events on a probability space  $S$ . Then

1.  $P(A^c) = 1 - P(A)$
2.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

**Proof.** For (1), we can use axioms 2 and 3 to write

$$1 = P(S) = P(A \cup A^c) = P(A) + P(A^c),$$

which is then equivalent to the claim. For (2), let's first notice (using Lemma 2.1(1)) that

$$\begin{aligned} P(A) &= P(A \cap S) = P(A \cap (B \cup B^c)) \\ &= P((A \cap B) \cup (A \cap B^c)) = P(A \cap B) + P(A \cap B^c). \end{aligned}$$

Also,

$$P(A \cup B) = P((A \cap B^c) \cup B) = P(A \cap B^c) + P(B),$$

from which we see that

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

Combining these relations leads to

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

which is equivalent to the claim. □

---

<sup>2</sup>Another apropos quote here is from the (Welsh-born) English philosopher and mathematician Bertrand Russell (1872–1970), who wrote, “The axiomatic method has many advantages over honest work.”

### 3 Counting Arguments

In our example in which a fair coin is flipped twice, we can compute probabilities simply by counting. To determine, for instance, the probability of getting both a head and a tail, we count the number of ways in which a head and a tail can both occur (2), and divide by the total number of possible outcomes (4). The probability is  $1/2$ . For large sample spaces, counting the number of possible outcomes in an event can be difficult, so it's useful to introduce some efficient techniques for proceeding. In these notes, we consider *permutations* and *combinations*, both of which can be understood in terms of the following simple rule for counting.

**Basic Principle of Counting.** *If  $N$  experiments are to be carried out, and there are  $n_1$  possible outcomes for the first experiment,  $n_2$  possible outcomes for the second experiment and so on up to  $n_N$  possible outcomes for the final experiment, then altogether there are*

$$n_1 \cdot n_2 \cdots n_N$$

*possible outcomes for the set of  $N$  experiments.*

**Example 3.1.** The *board* in Texas Hold'em poker consists of five cards dealt in order by the dealer (the first three are dealt together, but placed in order). If we keep track of order, how many boards are possible?

There are  $n_1 = 52$  ways the first card can be dealt,  $n_2 = 51$  remaining ways the second card can be dealt and so on, so that by using the basic principle of counting, we can compute

$$52 \cdot 51 \cdot 50 \cdot 49 \cdot 48 = 311,875,200.$$

△

#### 3.1 Permutations

Consider the following question: how many different numbers can be made by rearranging the digits 1, 2, and 3. There are  $n_1 = 3$  ways to choose the first digit,  $n_2 = 2$  remaining ways to choose the second digit, and  $n_3 = 1$  remaining way to choose the third digit, so by the basic principle of counting there are a total of  $3 \cdot 2 \cdot 1 = 6$  different numbers, namely 123, 132, 213, 231, 312, and 321. We refer to each of these arrangements as a *permutation*. As a general rule, we have the following:

**Rule of Permutations.** *For  $n$  distinct objects, there will be  $n!$  (read:  $n$  factorial) permutations, where*

$$n! = n \cdot (n - 1) \cdot (n - 2) \cdots 2 \cdot 1.$$

(Of course the 1 can be omitted, but it provides symmetry and closure, so it's been included.) The rule of permutations follows immediately from the basic principle of counting, through the observation that there are  $n$  ways to choose the first item (in our example, 1, 2 or 3),  $n - 1$  ways to choose the second item (once the first item has been eliminated) and so on.

**Example 3.2.** How many different four-letter permutations can be made from the letters in *math*?

We have four different letters, so we can have  $4! = 24$  permutations. △

**Example 3.3.** A certain student has ten books to arrange on a shelf: four math books, three chemistry books, two biology books, and one physics book. If the student wants to arrange the books so that all books dealing with the same subject are together, how many arrangements are possible?

First, there are  $4!$  ways to arrange the four subjects. Then there are  $4!$  ways to arrange the math books,  $3!$  ways to arrange the chemistry books,  $2!$  ways to arrange the biology books, and  $1!$  ways to arrange the physics book. In total, the number of arrangements is

$$4! \cdot 4! \cdot 3! \cdot 2! = 6,912.$$

△

In the event that one or more objects is repeated, we can establish a similar rule. Consider, for example, the following question: how many different five-digit numbers can be made by rearranging the digits 4, 5, 5, 5 and 6. First, if we had five different objects, say 4,  $A$ ,  $B$ ,  $C$ , and 6, then we would already know the answer:  $5!$ . But for each arrangement of 4 and 6, this count includes  $3!$  ways in which  $A$ ,  $B$ , and  $C$  can be arranged. E.g.,  $4ABC6$  is counted as different from  $4ACB6$ , and there are four other combinations with 4 at the left and 6 at the right. We conclude that the value  $5!$  overcounts by a factor of  $3!$ , so the correct count should be

$$\frac{5!}{3!} = 5 \cdot 4 = 20.$$

The general rule is as follows.

**Permutations with repeated objects.** For  $n$  objects for which  $n_1$  are identical,  $n_2$  are identical, etc., with  $n_N$  also identical, the number of permutations is given by

$$\frac{n!}{n_1!n_2! \cdots n_N!}.$$

**Example 3.4.** How many eight-letter permutations can be made from the letters in *Reveille*?

We have eight letters with one repeated three times (e) and one repeated twice (l). This gives

$$\frac{8!}{3!2!} = 3360$$

permutations. △

## 3.2 Combinations

We are often interested in counting the number of subsets we can create from some set of objects. For example, we might ask how many groups of three letters could be selected from the five letters  $A$ ,  $B$ ,  $C$ ,  $D$ , and  $E$ . If order matters, we argue that there are five ways to select the first letter (we have five possible options), four ways to select to the second (once

the first letter has been chosen, we only have four remaining to choose from), and three ways to select the third. That is, the number of possible selections can be computed as

$$5 \cdot 4 \cdot 3 = 60.$$

We observe, however, that this assumes the order of selection matters; that is, that the combination ABC is different from the combination BCA. When we talk about combinations, we will assume that order of selection *does not* matter, so the calculation above overcounts. In order to count the number of un-ordered combinations, we determine the number of ways in which the calculation above overcounts. For example, how many combinations have we counted that contain the letters A, B, and C? This is a permutation problem—how many ways can we permute A, B, and C—and the answer is  $3! = 6$ . Of course, we are overcounting every other combination of three letters by the same amount, so the total number of combinations is really

$$\frac{5 \cdot 4 \cdot 3}{3 \cdot 2 \cdot 1} = \frac{60}{6} = 10.$$

**Rule of Combinations.** *In general, if we have  $n$  distinct objects and choose  $r$ , we have the number of combinations*

$$\frac{n(n-1)(n-2)\cdots(n-(r-1))}{r!} = \frac{n!}{r!(n-r)!}.$$

For notational convenience, we make the following definition, typically read  *$n$  choose  $r$* ,

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

**Example 3.5.** The Texas Megamillions lottery works as follows: five different numbers are chosen in the range 1 to 75, and a sixth “Mega Ball” is chosen in the range of 1 to 15. (Order of the first five numbers does not matter, and indeed the numbers are always reported in ascending order.) How many possible combinations are there?

Using our rule of combinations, we easily compute this as

$$\binom{75}{5} 15 = 258,890,850.$$

△

While we’re talking about the lottery, consider the following oddity. On February 20, 2004 the jackpot for the Texas Megamillions lottery was \$230,000,000. It was under different rules at the time, with 135,145,920 possible winning combinations, with a cost of \$1.00 per ticket. This meant that in theory a player could play all 135,145,920 numbers and be assured of coming out ahead. In fact, of course, this is a dubious plan, since if someone else happens to pick the number as well, the player will have to share his winnings. Not to mention the logistics of buying this many tickets. But still.

**Example 3.6.** Determine the probability of obtaining each of the following poker hands if five cards are dealt from a standard deck of 52 cards: (1) straight flush, (2) four of a kind



(quad), (3) full house<sup>3</sup>, (4) flush, (5) straight, (6) three of a kind (trip), (7) two pair<sup>4</sup>, (8) one pair, (9) high-card hand (i.e., none of the above).

First, for this discussion the *rank* of a card refers to its value 2, 3, . . . , A, while the *suit* refers to the card's designation as clubs, diamonds, hearts, or spaces. As a starting point, since there are 52 distinct cards in a standard deck the total number of possible five card hands is

$$\binom{52}{5} = 2,598,960.$$

In each of the following calculations, we compute the probability of getting the result and not getting anything better. E.g., when we compute the probability of getting a pair, it is the probability that we get a pair and do not get three-of-a-kind, two pair etc.

1. *Straight flush*. Straights can be categorized by highest card: there are precisely ten rank arrangements possible, with high cards 5 through A. For each of these ten rank arrangements there are four suit arrangements that will make the straight a straight flush. This means:

$$\text{Number of possible straight flushes} = 10 \cdot 4 = 40.$$

Consequently the probability of getting a straight flush is

$$P(\text{straight flush}) = \frac{40}{2598960} = .000015.$$

We observe that there are precisely 4 ways to get a royal flush, so add a zero to the decimal.<sup>5</sup>

2. *Four of a kind*. There are 13 ways to rank the quad (i.e., 13 possible things to have four of) and 48 ways to choose the fifth card, so we have:

$$\text{Number of possible quads} = 13 \cdot 48 = 624.$$

Consequently the probability of getting four of a kind is

$$P(\text{four of a kind}) = \frac{624}{2598960} = .000240.$$

3. *Full house*. For a full house there are 13 ways to rank the trip and  $\binom{4}{3} = 4$  ways to arrange it, then 12 ways to rank the pair and  $\binom{4}{2} = 6$  ways to arrange it. We have:

$$\text{Number of possible full houses} = (13 \cdot 4) \cdot (12 \cdot 6) = 3744.$$

Consequently the probability of getting a full house is

$$P(\text{full house}) = \frac{3744}{2598960} = .001441.$$

---

<sup>3</sup>A full house is sometimes referred to as a *boat*, but while the phrases four of a kind and three of a kind seem clunky enough to warrant an abbreviation, full house does not.

<sup>4</sup>The word pair can be pluralized as either pairs or pair, with (I would suggest) pairs preferred in most circumstances. In the case of two pair poker hands it's traditional to pluralize with pair.

<sup>5</sup>On average, a player will be dealt a royal flush once in every 649,740 hands.

4. *Flush (and not a straight flush)*. We have four ways to choose the suit for a flush and  $\binom{13}{5} = 1287$  ways to arrange the ranks. Finally, we will subtract off the 40 straight flushes we've already considered. We have

$$\text{Number of possible flushes (not straight flushes)} = 4 \cdot 1287 - 40 = 5,108.$$

Consequently the probability of getting a flush (and not a straight flush) is

$$P(\text{flush (and not a straight flush)}) = \frac{5108}{2598960} = .001965.$$

5. *Straight (and not a straight flush)*. As in our discussion of straight flushes, we first note that there are ten possible rank arrangements for a straight, corresponding with top cards 5 through A. For each such rank arrangement there are four ways to suit each of the five cards. Finally, we subtract off the straight flushes. This gives:

$$\text{Number of straights (not straight flushes)} = 10 \cdot 4^5 - 40 = 10,200.$$

Consequently the probability of getting a straight (and not a straight flush) is

$$P(\text{straight (and not a straight flush)}) = \frac{10200}{2598960} = .003925.$$

6. *Three of a kind*. We have 13 ways to choose the rank of the trip,  $\binom{4}{3} = 4$  ways to arrange the trip,  $\binom{12}{2} = 66$  ways to choose two new ranks (which must be different to avoid a full house), and 4 ways to suit each of these ranks. (Notice that we have no problem overcounting flushes or straights since no hand with a trip can be a flush or a straight.) We have, then

$$\text{Number of possible trips} = 13 \cdot 4 \cdot 66 \cdot 4^2 = 54,912.$$

Consequently the probability of getting three of a kind is

$$P(\text{three of a kind}) = \frac{54912}{2598960} = .021128.$$

7. *Two pair*. We have  $\binom{13}{2} = 78$  ways to choose the two different ranks to be paired,  $\binom{4}{2} = 6$  ways to arrange each pair, and 44 ways to pick the remaining card. In total,

$$\text{Number of possible two-pair hands} = 78 \cdot 6^2 \cdot 44 = 123,552.$$

Consequently the probability of getting two pair is

$$P(\text{two-pair hand}) = \frac{123552}{2598960} = .047539.$$

8. *One pair*. We have 13 ways to rank the pair and  $\binom{4}{2} = 6$  ways to arrange it. There are  $\binom{12}{3} = 220$  ways to rank the remaining three cards so that no two of the ranks are the same, and four ways to suit each. We have:

$$\text{Number of possible one-pair hands} = 13 \cdot 6 \cdot 220 \cdot 4^3 = 1,098,240.$$

Consequently the probability of getting one pair is

$$P(\text{one-pair}) = \frac{1098240}{2598960} = .422569.$$

9. *High-card hand.* Summing all hands considered in (1) through (8) we have 1,296,420 hands. This leaves 1,302,540 hands that are ranked only by the ranks of the individual cards. We refer to these as high card hands.<sup>6</sup> We have:

$$P(\text{high-card hand}) = \frac{1302540}{2598960} = .501177.$$

Alternatively, we can check our calculations by computing the number of high-card hands directly. There are  $\binom{13}{5} = 1287$  ways to choose five distinct ranks and four ways to suit each of the ranks. We must subtract flushes, straights, and straight flushes, all of which are created with five distinct ranks. That is:

$$\text{Number high card hands} = 1287 \cdot 4^5 - 40 - 5108 - 10200 = 1302540,$$

as above. △

**Example 3.7.** (Warm-up for the next topic) Suppose we have three stars and four bars that we would like to arrange in a line such as

$$** || * || \quad \text{or} \quad | * | * | * |$$

How many arrangements are possible?

We can view this as having seven slots to fill, and we can choose three to fill with a star. I.e.,

$$\binom{7}{3} = 35.$$

Of course, we could make exactly the same argument focusing on bars, in which case we would compute

$$\binom{7}{4} = 35.$$

### 3.3 Combinations with Repetition

Now, suppose we want to know how many three-letter combinations we can make from the letters  $A$ ,  $B$ ,  $C$ ,  $D$ , and  $E$ , but in this case we will allow letters to repeat, so e.g.,  $AAB$  would be valid, or  $CCC$  etc. Here's a standard way of thinking about this: let  $x_1$  denote the number of times  $A$  appears,  $x_2$  the number of times  $B$  appears, and so on, up to  $x_5$  being the number of times  $E$  appears. Then

$$x_1 + x_2 + x_3 + x_4 + x_5 = 3.$$

In fact, the number of non-negative solutions to this equation is exactly what we're looking for. We can approach this by the method of "stars and bars." For example, consider the

---

<sup>6</sup>A.k.a., bupkis.

solution  $x_1 = 2$ ,  $x_3 = 1$ , with the rest 0. We denote the numbers  $x_i$  by stars, and divide them by bars:

$$\underbrace{\star\star}_{x_1} | \underbrace{\quad}_{x_2} | \underbrace{\star}_{x_3} | \underbrace{\quad}_{x_4} | \underbrace{\quad}_{x_5}$$

Without the extra baggage, the stars and bars arrangement is just

$$\star\star || \star ||$$

But we already know how to compute the number of such arrangements,

$$\binom{7}{3} = 35.$$

**Combinations with repetition.** *In general, if we have  $n$  distinct objects divided into groups of  $r$  with repeats allowed, then we have  $r$  stars and  $n - 1$  bars so*

$$\binom{n+r-1}{r}$$

*possible combinations.*

### 3.4 Multiple Combinations

Our next question regards the following situation: suppose a set of  $n$  distinct objects is to be divided into  $N$  distinct groups of sizes  $n_1, n_2, \dots, n_N$ , where

$$n_1 + n_2 + \dots + n_N = n.$$

Notice that a single combination is the case  $N = 2$ , with  $n_1 = r$  and  $n_2 = n - r$ .

To be specific, let's suppose  $n = 10$  students are to be assigned to  $N = 3$  different projects,  $A$ ,  $B$ , and  $C$ , with  $n_1 = 3$  students assigned to project  $A$ ,  $n_2 = 3$  students assigned to project  $B$ , and  $n_3 = 4$  students assigned to project  $C$ . The question is: in how many ways can the students be distributed to the projects? In this case, there are  $\binom{10}{3}$  ways to assign the first group,  $\binom{7}{3}$  remaining ways to assign the second group, and  $\binom{4}{4}$  remaining ways to assign the third. In total, we have

$$\binom{10}{3} \binom{7}{3} \binom{4}{4} = 4,200.$$

Notice that the general case is now becoming clear. There will be  $\binom{n}{n_1}$  ways to arrange the first group,  $\binom{n-n_1}{n_2}$  ways to arrange the second group, and so on, until

$$\binom{n - n_1 - n_2 - \dots - n_{N-1}}{n_N}$$

ways to arrange the  $N^{\text{th}}$  group. The total number of multiple combinations is

$$\begin{aligned} & \binom{n}{n_1} \binom{n-n_1}{n_2} \cdots \binom{n-n_1-n_2-\cdots-n_{N-1}}{n_N} \\ &= \frac{n!}{(n-n_1)!n_1!} \cdot \frac{(n-n_1)!}{(n-n_1-n_2)!n_2!} \cdots \frac{(n-n_1-\cdots-n_{N-1})}{(n-n_1-\cdots-n_N)n_N!} \\ &= \frac{n!}{n_1!n_2!\cdots n_N!}. \end{aligned}$$

Notice that this is exactly the same formula we found when looking at permutations with repeats. Why? For our example with 10 students and three groups, we could think about having the letters  $AAABBBCCCC$ , and distributing these letters among some ordering of the students, e.g., Rose gets the first letter, Colin gets the second letter etc. In this way, each permutation of the letters corresponds with a unique assignment of groups, so the two counts are indeed identical.

In light of this, we generalize our “choose” notation by writing

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

where it is assumed that

$$n_1 + n_2 + \cdots + n_N = n.$$

These values are called *multinomial coefficients* because of their appearance in the *multinomial formula*

$$(x_1 + x_2 + \cdots + x_N)^n = \sum_{n_1+\cdots+n_N=n} \binom{n}{n_1, n_2, \dots, n_N} x_1^{n_1} x_2^{n_2} \cdots x_N^{n_N}.$$

**Side note.** To get a feel for what the multinomial formula tells us, let’s consider the case with  $N = 2$  and  $n = 3$ , for which we get

$$(x_1 + x_2)^3 = \sum_{k=0}^3 \binom{3}{k} x_1^k x_2^{3-k}.$$

Of course, we can also write

$$(x_1 + x_2)^3 = (x_1 + x_2)(x_1 + x_2)(x_1 + x_2).$$

If we think about multiplying out the right-hand side of this last expression, we can think of the coefficient  $\binom{3}{k}$  as counting the number of ways  $k$  we can choose  $x_1$  from each of the three factors. E.g.,  $k = 0$  corresponds with choosing  $x_1$  zero times, leaving  $x_2^3$ . There is only one way to do this (i.e.,  $\binom{3}{0} = 1$ ). Likewise,  $k = 1$  corresponds with choosing  $x_1$  once, and there are  $\binom{3}{1} = 3$  ways to do this (once from each factor). The cases  $k = 2$  and  $k = 3$  are similar.

Returning to our example with 10 students, let’s make things slightly more complicated. Previously, we assumed it was particularly project  $C$  that would have four students, but

suppose instead that the group of four students could be assigned to any of the three projects, with still three students assigned to the other two. In this case, we start by observing that there are three ways to assign the group of four. Once that's been assigned, there aren't really two ways to assign the first group of three, because the order in which we assign the 3-student projects doesn't matter. I.e., if we multiplied by 2, we would be counting the same set of groups twice, once if we assigned the students first to  $A$ , and once if we assigned the students first to  $B$ . This means that the actual number of assignments in this case is

$$3 \cdot 4,200 = 12,600.$$

Finally, suppose the projects are all identical, and we only want to know the number of possible student groupings. In this case, our original calculation double-counted each grouping, because it distinguished between whether a student was assigned to project  $A$  or to project  $B$ . To be clear about this, let's denote the students  $\{S_i\}_{i=1}^{10}$ , and think about one way in which we might arrange them:

$$\begin{aligned} A &: S_1 S_2 S_3 \\ B &: S_4 S_5 S_6 \\ C &: S_7 S_8 S_9 S_{10} \end{aligned}$$

For our first count, this would differ from the arrangement

$$\begin{aligned} A &: S_4 S_5 S_6 \\ B &: S_1 S_2 S_3 \\ C &: S_7 S_8 S_9 S_{10} \end{aligned}$$

while for the current count it does not. So, for the current count, we end up with

$$\frac{4,200}{2} = 2,100$$

groupings.

## 4 Conditional Probability

Often, we would like to compute the probability that some event occurs, given a certain amount of information. In our two-flip game above, suppose that we are given the information that the first flip is heads (i.e., event  $E$ ), and we would like to compute the probability of getting two heads. We would write

$$P((HH)|E) = \frac{1}{2},$$

and read this as the probability of getting (HH) "given" event  $E$ . Alternatively, suppose we are given the information that heads turns up on at least one flip,

$$G = \{(HH), (HT), (TH)\}.$$

Then

$$P((HH)|G) = \frac{1}{3}.$$

Notice that conditioning has the effect of reducing the size of the sample space (to  $E$  in the first case, and to  $G$  in the second).

**Definition 4.1.** (*Conditional probability*) Suppose  $A$  and  $B$  are events on a sample space  $S$  and  $P(B) \neq 0$ . Then we define the conditional probability of event  $A$  given that event  $B$  has occurred as

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

**Justification for the definition.** First, let's check that this definition agrees with our above intuitive calculation. For the first, we would write

$$P((HH)|E) = \frac{P((HH) \cap E)}{P(E)} = \frac{P(HH)}{P(E)} = \frac{\frac{1}{4}}{\frac{1}{2}} = \frac{1}{2},$$

while for the second we would write

$$P((HH)|G) = \frac{P((HH) \cap G)}{P(G)} = \frac{P(HH)}{P(G)} = \frac{\frac{1}{4}}{\frac{3}{4}} = \frac{1}{3}.$$

More generally, we recall that our intuitive way of thinking about probabilities is to imagine carrying out experiments and computing

$$P(A|B) \cong \frac{\# \text{ of times } A \text{ occurs along with } B}{\# \text{ of times } B \text{ occurs}}.$$

I.e., we begin by conducting experiments, and if  $B$  does not occur then we discard the result. If  $B$  does occur, then we check if  $A$  occurred as well, and if so add a count to the numerator. If  $n$  denotes the number of experiments, then we can write

$$P(A|B) \cong \frac{\frac{\# \text{ of times } A \text{ occurs along with } B}{n}}{\frac{\# \text{ of times } B \text{ occurs}}{n}}.$$

Here, the numerator is precisely how we approximate  $P(A \cap B)$ , and the denominator is how we approximate  $P(B)$ . In the limit as  $n \rightarrow \infty$ , we (intuitively) recover the definition of conditional probability.

**Example 4.1.** Suppose two fair dice are rolled. What is the probability that at least one lands on six, given that the dice land on different numbers?

Let  $A$  be the event of at least one die landing on six, and let  $B$  be the event that the dice land on different numbers. We immediately see by counting outcomes that  $P(B) = \frac{30}{36} = \frac{5}{6}$ . On the other hand, the probability of 1 six and 1 non-six is the number of possible combinations with exactly one six (10) divided by the total number of possible combinations (36):

$$P(A \cap B) = P(1 \text{ six}, 1 \text{ not-six}) = \frac{\# \text{ combinations with 1 six, 1 not-six}}{36 \text{ total combination possible}} = \frac{10}{36} = \frac{5}{18}.$$

Consequently,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{\frac{5}{18}}{\frac{5}{6}} = \frac{1}{3}.$$

△

## 4.1 Computing Probabilities by Conditioning

In many cases, it's useful to turn the definition of conditional probability around and write

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

and more generally

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

Continuing in this way (and turning things around to get a convenient form), we arrive at the general relation for sets  $\{A_k\}_{k=1}^n$ ,

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

**Example 4.2.** In a standard game of bridge, what is the probability that each player will hold exactly one ace? Here, the only thing we need to know about bridge is that it's played with a standard deck of 52 cards, and each of four players receives 13 cards.

We set

$A_1$  = event the aces of clubs and diamonds are in two different hands

$A_2$  = event the aces of clubs, diamonds, and hearts are in three different hands

$A_3$  = event the four aces are in four different hands.

Noticing that  $A_3 \subset A_2 \subset A_1$ , our set-up will be

$$P(A_3) = P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2).$$

For  $P(A_1)$ , the ace of clubs must land somewhere, so all we have to compute is the probability that the ace of diamonds doesn't land the same place. For the three hands that do not contain the ace of clubs, 39 of 51 cards will be dealt, so

$$P(A_1) = \frac{39}{51}.$$

Just to be clear about this, notice that there will be  $\binom{51}{39}$  ways to fill the hands that do not contain the ace of clubs, and if we want to ensure that the ace of diamonds appears in one of these hands then we are left with 50 cards to distribute to the other 38 places. In total,

$$P(A_1) = \frac{\binom{50}{38}}{\binom{51}{39}} = \frac{\frac{50!}{12!38!}}{\frac{51!}{12!39!}} = \frac{39}{51}.$$



For  $P(A_2|A_1)$ , we proceed similarly, noting that the ace of hearts must now land in one of two remaining hands, with 26 cards out of 50 to be dealt. I.e.,

$$P(A_2|A_1) = \frac{26}{50} = \frac{13}{25}.$$

Last, and by similar reasoning,

$$P(A_3|A_1 \cap A_2) = \frac{13}{49}.$$

In total,

$$P(A_3) = \frac{39}{51} \cdot \frac{13}{25} \cdot \frac{13}{49} = .1055.$$

△

**Example 4.3.** Suppose there are 20 students in a certain modeling class. What is the probability that at least two of the students share the same birthday (day, not necessarily year)? We will ignore leap years and century leap years, though these can be taken into account similarly with more work.

By assumption, we are taking the number of days in the year to be precisely 365. We will think of labeling the students as student 1, student 2, and so on, and we will specify the events

- $A_1$  = event students 1 and 2 have different birthdays
- $A_2$  = event students 1, 2, and 3 have different birthdays
- $\vdots$
- $A_{19}$  = event all 20 students have different birthdays.

Our goal is to compute

$$P(A_{19}^c) = 1 - P(A_{19}).$$

Similarly as in the previous example, we will proceed by computing

$$\begin{aligned} P(A_{19}) &= P(A_1 \cap A_2 \cap \cdots \cap A_{19}) \\ &= P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2) \cdots P(A_n|A_1 \cap A_2 \cap \cdots \cap A_{18}). \end{aligned}$$

For  $P(A_1)$ , student 1 has some birthday, so we only need to compute the probability that student 2 has a different birthday. For this, there are 364 out of 365 days available, so  $P(A_1) = \frac{364}{365}$ . Likewise,  $P(A_2|A_1) = \frac{363}{365}$ , where in this case the denominator does not change, because using a day of the year doesn't reduce the number of days possible. Continuing in this way, the final probability we need is

$$P(A_n|A_1 \cap A_2 \cap \cdots \cap A_{18}) = \frac{346}{365},$$

and we can compute

$$P(A_{19}) = \frac{364}{365} \cdot \frac{363}{365} \cdots \frac{346}{365} = .5886.$$

I.e., the probability we're looking for is

$$P(A_{19}^c) = 1 - P(A_{19}) = 1 - .5886 = .4114.$$

△

**Definition 4.2.** We say that the events  $\{A_k\}_{k=1}^n$  partition a sample space  $S$  if they are mutually exclusive, and

$$\bigcup_{k=1}^n A_k = S.$$

For example, the collection of individual outcomes in a sample space always forms a partition of  $S$ , as does the space  $S$  itself. For our two-flip experiment, one partition of  $S$  would be  $A_1 = \{(HH), (HT)\}$  and  $A_2 = \{(TH), (TT)\}$ .

Given a partition of  $S$ ,  $\{A_k\}_{k=1}^n$  and any additional event  $B \subset S$ , we can write

$$B = B \cap S = B \cap \left(\bigcup_{k=1}^n A_k\right) = \bigcup_{k=1}^n (B \cap A_k),$$

where the final equality follows from a straightforward extension of Lemma 2.1(1). The collection  $\{B \cap A_k\}_{k=1}^n$  is mutually exclusive, so we find

$$P(B) = P\left(\bigcup_{k=1}^n (B \cap A_k)\right) = \sum_{k=1}^n P(B \cap A_k) = \sum_{k=1}^n P(B|A_k)P(A_k).$$

This gives us another useful way for computing probabilities by conditioning.

**Example 4.4.** Suppose a certain baseball player can hit a fast ball 10% of the time, a knuckle ball 20% of the time, and a curve ball 30% of the time, and he's up against a pitcher who throws these pitches respectively 40%, 40%, and 20% of the time. What is the probability the player will get a hit on any particular pitch?

We define the events

$B$  = event player gets a hit

$A_1$  = event of a fast ball

$A_2$  = event of a knuckle ball

$A_3$  = event of a curve ball

Then

$$\begin{aligned} P(B) &= P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3) \\ &= .1 \cdot .4 + .2 \cdot .4 + .3 \cdot .2 = .18. \end{aligned}$$

△

**Example 4.5.** Suppose a fair pair of dice is rolled repeatedly, and the sum of values is recorded each time. What is the probability that a sum of 5 will occur before a sum of 7?

For this one, we specify the following events:

$B$  = event a sum of 5 occurs before a sum of 7

$A_1$  = event a sum of 5 occurs on the first roll

$A_2$  = event a sum of 7 occurs on the first roll

$A_3$  = event neither a sum of 5 nor a sum of 7 occurs on the first roll

Here,  $A_1$  occurs with a roll of 1 and 4 or a roll of 2 and 3, so

$$P(A_1) = \frac{4}{36} = \frac{1}{9}.$$

Likewise,  $A_2$  occurs with a rolls of 1 and 6, 2 and 5, or 3 and 4, giving

$$P(A_2) = \frac{6}{36} = \frac{1}{6}.$$

For  $A_3$ , we can compute

$$P(A_3) = 1 - P(A_1) - P(A_2) = 1 - \frac{1}{9} - \frac{1}{6} = \frac{13}{18}.$$

This allows us to compute

$$\begin{aligned} P(B) &= P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3) \\ &= 1 \cdot \frac{1}{9} + 0 \cdot \frac{1}{6} + P(B) \cdot \frac{13}{18}. \end{aligned}$$

This is an equation that we can solve for  $P(B)$ , and we find that

$$P(B) = \frac{\frac{1}{9}}{1 - \frac{13}{18}} = \frac{\frac{1}{9}}{\frac{5}{18}} = \frac{2}{5}.$$

△

## 4.2 Bayes' Lemma

The next result we will develop is called Bayes' Lemma (often Bayes' Theorem), named for the British statistician Thomas Bayes (1701–1761).

**Lemma 4.1.** (*Bayes' Lemma*) Suppose the events  $\{A_k\}_{k=1}^n$  form a partition of a sample space  $S$ . If  $B$  is an event in  $S$  and  $P(B) \neq 0$ , then for each  $k \in \{1, 2, \dots, n\}$  we have,

$$P(A_k|B) = \frac{P(B|A_k)P(A_k)}{\sum_{j=1}^n P(B|A_j)P(A_j)}.$$

**Proof.** Using our previous calculations, we can write

$$P(A_k|B) = \frac{P(A_k \cap B)}{P(B)} = \frac{P(B|A_k)P(A_k)}{\sum_{j=1}^n P(B|A_j)P(A_j)},$$

which is precisely the claim. □

**Example 4.6.** Suppose a laboratory test is 95% effective in detecting a certain disease when the disease is present, but the test also yields a fake positive result for 1% of healthy people tested. If .5% (i.e., .005) of the population actually has the disease, what is the probability that a person has the disease, given that the test result is positive?

For this example, we specify the events

$A$  = event a tested person has the disease

$B$  = event of a positive test

The probability we want to compute is  $P(A|B)$ . In order to use Bayes' Lemma, we introduce the simple partition  $A_1 = A$  and  $A_2 = A^c$ . This allows us to compute

$$\begin{aligned} P(A|B) &= P(A_1|B) = \frac{P(B|A_1)P(A_1)}{P(B|A_1)P(A_1) + P(B|A_2)P(A_2)} \\ &= \frac{.95 \cdot .005}{.95 \cdot .005 + .01 \cdot .995} = .3231. \end{aligned}$$

I.e., if the test comes back positive, there is only about a 32% chance the person actually has the disease. △

**Example 4.7.** (The infamous Monty Hall problem.<sup>7</sup>) Consider a game show in which a prize is hidden behind one of three doors. The contestant chooses a door, and then the host opens one of the two unchosen doors, showing that the prize is not behind it. (He never opens the door that the prize is behind.) The contestant then gets the option to switch doors. Given this scenario, should a contestant hoping to optimize his winnings, (1) always switch doors, (2) never switch doors, or (3) doesn't matter?

Though we can argue a solution to this problem on intuitive grounds (be warned: the argument might not be the first one you think of), we will work through the details as an application of Bayes' Lemma. We will determine the probability that the prize is behind the first door the contestant selects, given that the host opens one of the other doors. We begin by defining a set of appropriate events, in which the doors the contestant does not open are generically labeled *alternative door number 1* and *alternative door number 2*:

$A_1$  = event that prize is behind first door selected

$A_2$  = event that prize is behind alternative door 1

$A_3$  = event that prize is behind alternative door 2

$B$  = event that host opens alternative door 1

---

<sup>7</sup>Monty Hall was the host of a show called *Let's Make a Deal*, which ran from 1963 to 1976. This problem was first proposed by the American biostatistician Steve Selvin (born 1941) in 1975. The problem also appeared in the September 1990 version of Marilyn vos Savant's long-running column *Ask Marilyn*. Marilyn got the answer correct, but received over 10,000 letters, many from mathematics and statistics professors, claiming she was wrong. The problem was also popularized by the 2008 Robert Luketic film *21*. For an interactive version of the game that you can play on-line, see <http://www.shodor.org/interactivate/activities/monty3/>

According to Bayes' Lemma, we have

$$\begin{aligned} P(A_1|B) &= \frac{P(B|A_1)P(A_1)}{P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3)} \\ &= \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2} \cdot \frac{1}{3} + 0 + 1 \cdot \frac{1}{3}} = \frac{\frac{1}{6}}{\frac{1}{2}} = \frac{1}{3}. \end{aligned}$$

Observe that since the host's opening alternative door 1 is entirely arbitrary, we can make the same calculation given that he opens alternative door 2. Therefore, whichever door he opens, the contestant who sticks with his initial choice only has a  $1/3$  chance of being right. Which, of course, means that the contestant should always switch doors, giving him a  $2/3$  chance of winning.

Intuitively, we can regard this game as follows: If the contestant chooses not to switch doors, his only chance of winning is if his original choice was correct, and his odds for that are clearly  $1/3$ . On the other hand, since the host removes one of the two doors not selected, if the contestant switches doors he wins so long as his original choice was incorrect.  $\triangle$

**Remark.** *In many applications of Bayes' Lemma, it is not immediately clear that the conditioning set  $B$  is in the sample space  $S$ . For example, in the Monty Hall problem, the events that partition  $S$  correspond only with possible locations of the prize and seem to have nothing to do with which door the host opens. In such cases, it is important to keep in mind that while a set of outcomes entirely describes  $S$ , a set of events can hide individual outcomes. In order to see this more clearly, notice that*

$$B^c = \text{event the host opens alternative 2.}$$

Then we can write the sample space for the Monty Hall problem as follows:

$$S = \{(A_1 \cap B), (A_1 \cap B^c), (A_2 \cap B^c), (A_3 \cap B)\}.$$

Here  $A_1 = \{(A_1 \cap B), (A_1 \cap B^c)\}$ ,  $A_2 = \{(A_2 \cap B^c)\}$ ,  $A_3 = \{(A_3 \cap B)\}$ , and so the  $A_k$  clearly partition  $S$ . Moreover,  $B = \{(A_1 \cap B), (A_3 \cap B)\}$ , which is clearly a subset of  $S$ .

### 4.3 Independent Events

If the occurrence of an event  $B$  gives no information about an event  $A$ , then

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

For example, in our two-flip experiment suppose  $A = \{(HH), (HT)\}$  is the event that heads occurs on the first flip and  $B = \{(HT), (TT)\}$  is the event that tails occurs on the second flip. Then

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

This agrees with our intuition that the appearance of tails on the second flip gives no information about whether heads occurred on the first flip. Notice that if  $P(A|B) = P(A)$ , then

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

and from this we see that

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

giving an expected symmetry. Precisely, we make the following definition.

**Definition 4.3.** We say that two events  $A$  and  $B$  are independent provided

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

If two events are not independent, then they are said to be dependent.

We typically use this latter characterization of independence, because it's valid in the case  $P(B) = 0$ . Otherwise, it's equivalent to the statement  $P(A|B) = P(A)$ .

**Example 4.8.** Compute the probability that double sixes occurs at least once in 24 rolls of a fair pair of dice.

We begin this calculation by computing the probability that double sixes does not occur even once in the 24 throws. On each trial the probability that double sixes will not occur is  $\frac{35}{36}$ , and so by independence the probability of the intersection of 24 of these events in a row is  $(\frac{35}{36})^{24}$ . We conclude that the probability that double sixes *do* turn up at least once in 24 throws is (to four decimal places of accuracy)

$$1 - \left(\frac{35}{36}\right)^{24} = .4914.$$

△

## 5 Discrete Random Variables

Suppose that in our experiment of flipping a coin twice, we assigned a numerical value to each outcome, referred to as  $X$ :  $X(HH) = 1$ ,  $X(HT) = 2$ ,  $X(TH) = 3$ , and  $X(TT) = 4$ . For instance, we might be considering a game in which  $X$  represents the payoff for each possible outcome. (Below, we will refer to this game as the “two-flip game.”) We refer to functions defined on sample spaces as *random variables*, and we refer to the values random variables can take as *realizations*. Random variables represent such processes as:

- The value of a stock at a given time;
- The amount of money a player wins in a game of chance;
- The time it takes to check out at a supermarket.

Random variables that can take only a countable number of values are called *discrete*. (Recall that a set is said to be *countable* if its elements can be enumerated 1, 2, 3, .... The set of all rational numbers (integer fractions) is countable; the set of all real numbers is not.)

We will define events with respect to random variables in the forms  $\{X = 1\}$ ,  $\{X \leq 3\}$ ,  $\{X \geq 2\}$  etc., by which we mean the event in our sample space for which  $X$  satisfies the condition in brackets. For example,  $\{X = 1\} = \{(HH)\}$  and  $\{X \leq 2\} = \{(HH), (HT)\}$ .

## 5.1 Expected Value

Often, we would like to summarize information about a particular random variable. For example, we might ask how much we could expect to make playing the two-flip game. Put another way, we might ask, how much would we make on average if we played this game repeatedly a sufficient number of times. In order to compute this *expected value*, we multiply the amount we win from each outcome with its probability and sum. In the case of the two-flip game, we have

$$\text{Expected value} = \$1.00 \times \frac{1}{4} + \$2.00 \times \frac{1}{4} + \$3.00 \times \frac{1}{4} + \$4.00 \times \frac{1}{4} = \$2.50.$$

It's important to notice that we will never actually make \$2.50 in any single play of the game. But if we play it repeatedly for a sufficient length of time, our average winnings will be \$2.50 per game played. Denoting expected value by  $E$ , we summarize this critical expression as

$$E[X] = \sum_{x \text{ Possible}} xP(X = x).$$

**Example 5.1.** Suppose a man counting cards at the blackjack table knows the only cards not yet dealt are a pair of fours, three nines, a ten, two Jacks, and a King. What is the expected value of his next card?

Keeping in mind that tens, Jacks and Kings are all worth ten points, while fours and nines are worth face value, we compute

$$E[\text{next card drawn}] = 4 \cdot \frac{2}{9} + 9 \cdot \frac{3}{9} + 10 \cdot \frac{4}{9} = \frac{25}{3}.$$

△

**Example 5.2.** The game of (American) roulette employs a large wheel with 38 slots, thirty-six labeled nonsequentially by the numbers 1–36, and two house slots, labeled 0 and 00.<sup>8</sup> Exactly half of the thirty-six numbered slots are red, while the other half are black. (The two house slots are typically shaded green.) In each play of the game, players make bets on certain numbers or sequences of numbers until the *croupier*<sup>9</sup> calls out “No more bets,” usually while the wheel is already in motion, but before it would be possible for anyone to guess where the ball he’s spun into it might land. The simplest bet in roulette is on red or black, which returns even money (the player betting one dollar either loses his dollar or wins one dollar). The expected value of a one dollar bet on black can easily be computed as,

$$E[\text{one dollar bet on black}] = 1 \cdot \frac{18}{38} - 1 \cdot \frac{20}{38} = -.0526 \text{ cents},$$

with variance

$$\text{Var}[\text{one dollar bet on black}] = 1.0526^2 \cdot \frac{18}{38} + .9474^2 \cdot \frac{20}{38} = .9972.$$

---

<sup>8</sup>In European roulette, there is only one house slot, 0.

<sup>9</sup>French for “dealer.”

On the other extreme, we could bet on a single number, which pays 35:1 (\$35 dollar won for \$1 bet) giving an expected value

$$E[\text{Single Number}] = 35 \cdot \frac{1}{38} - 1 \cdot \frac{37}{38} = -.0526.$$

One well-known gambling strategy that is often used in roulette, though is applicable to just about any game, is referred to as the “double-up”, or *Martingale* strategy. The idea behind the Martingale strategy is extremely simple: At every play of the game, if the player wins, he collects his money and repeats his initial bet, starting over; if he loses, he doubles his bet *on the same thing*. For example, a player might use this strategy betting on black in roulette. He begins by betting, say, one dollar on black. If he wins, he collects the dollar he wins and starts over, betting a single dollar on black again. If he loses, he bets two dollars on black. Now, if he wins, he collects two dollars from the croupier, which even counting the bill he’s already lost leaves him ahead by \$1. If he loses, he doubles his bet again to four dollars. Assuming black eventually turns up, then no matter how long it takes for him to get there, the player will eventually make \$1. He can’t lose, right? Well...unfortunately, there’s some fine print. True, the mathematical expectation of this strategy is \$1—you really do, in theory, make a buck every time you use it. The problem is that casinos have betting limits (and gamblers have capital limits), so eventually the doubling has to stop. In order to see what effect betting limits have on the expectation, let’s suppose there’s an extremely small cap of \$4. If we let  $X$  be the player’s final payout from one series with the strategy, we have

$$E[X] = 1 \cdot \left( \frac{18}{38} + \frac{20}{38} \cdot \frac{18}{38} + \left( \frac{20}{38} \right)^2 \frac{18}{38} \right) - 7 \left( \frac{20}{38} \right)^3 = -.1664.$$

Here, the values  $\frac{18}{38}$ ,  $\frac{20}{38} \cdot \frac{18}{38}$ , and  $\left( \frac{20}{38} \right)^2 \frac{18}{38}$  are the respectively the probabilities that the player wins on the first, second, and third spins (winning \$1.00 in each case), and  $\left( \frac{20}{38} \right)^3$  is the probability that she loses on each of the first three spins.

At first glance, it might appear that the player’s expected value is getting worse, but the situation isn’t quite that bad. Notice that the player’s bet is no longer a single dollar, but rather varies depending upon the turn of the wheel. His expected bet under the assumption of a \$4 cap is given by,

$$E[B] = 1 \cdot \frac{18}{38} + 3 \cdot \frac{20}{38} \cdot \frac{18}{38} + 7 \cdot \left( 1 - \frac{18}{38} - \frac{20}{38} \cdot \frac{18}{38} \right) = 3.1607.$$

(Notice that in this calculation the first addend on the right-hand side represents the player’s bet in the event that he bets exactly one dollar, which only occurs in the event that he wins on the first roll, and the others are similar.) The player’s expectation *per dollar bet* becomes

$$\text{Expectation per dollar bet} = \frac{E[X]}{E[B]} = \frac{-.1664}{3.1607} = -.0526.$$

So, in fact, the player hasn’t helped (or hurt) herself at all. △

**Example 5.3.** (St. Petersburg Paradox, suggested by Daniel and Nicolaus Bernoulli around 1725.) Suppose a dealer says that he will flip a fair coin until it turns up heads and will pay



the player  $2^n$  dollars, where  $n$  is the number of flips it takes for the coin to land heads. How much should a player be willing to pay in order to play this game?

When asked this question, most people say that the player shouldn't pay much over \$2.00, and just about everyone agrees that she shouldn't go above \$5.00. In order to determine the expected value of this game, let's let  $X$  denote a random variable giving the game's the payoff on a given play. If the coin lands heads on the first flip the payoff is \$2, with probability  $\frac{1}{2}$ . If the coin does not land heads until the second flip, the payoff is  $2^2 = 4$ , with probability  $\frac{1}{4}$ —the probability of a tail followed by a head. Proceeding similarly, we see that

$$E[X] = 2 \cdot \frac{1}{2} + 4 \cdot \frac{1}{4} + 8 \cdot \frac{1}{8} + \dots = 1 + 1 + 1 + \dots = \infty.$$

The expected value of this game is infinite! Which means that *ideally* we should be willing to pay any amount of money to play it. But almost no one is willing to pay more than about five bucks. This is what the brothers Bernoulli considered a paradox.

In order to resolve this, we need to keep in mind that the expected value of a game reflects a player's average winnings if she were to continue playing the game for a sufficient length of time. Suppose she pays five dollars per game. Half the time she will lose three dollars ( $2^1 - 5 = -3$ ), while another quarter of the time she will lose one dollar ( $2^2 - 5 = -1$ ). On the other hand, roughly one out of every sixty-four times (6 flips) she will make  $2^6 - 5 = 59$ . The point is that though the player will lose more often than she will win, she will occasionally win big. Practically speaking, this means that two things come into play when thinking about the expected value of a game: the expected value itself and the number of times the player will get a chance to play it. Yet one more way to view this is as follows. The fact that this game has infinite expectation means that no matter how much the dealer charges someone to play—\$5.00, \$500.00, \$5 million, etc.—the game is worth playing (i.e., the player will eventually come out ahead) *as long as the player can be sure that she will be able to play it enough times.*  $\triangle$

**Example 5.4.** Compute the expected value for a *binomial* random variable with sample size  $n$  and probability  $p$ , which is a random variable  $X$  that takes values  $0, 1, 2, \dots, n$  with probability

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

(See Section 5.4 for additional information about binomial random variables.)

We compute

$$\begin{aligned} E[X] &= \sum_{k=0}^n k p(k) = \sum_{k=1}^n k \binom{n}{k} p^k (1 - p)^{n-k} = \sum_{k=1}^n k \frac{n!}{k!(n-k)!} p^k (1 - p)^{n-k} \\ &= \sum_{k=1}^n \frac{n!}{(k-1)!(n-k)!} p^k (1 - p)^{n-k}. \end{aligned}$$

Letting  $l = k - 1$ , we have

$$\begin{aligned} E[X] &= \sum_{l=0}^{n-1} n \frac{(n-1)!}{l!((n-1)-l)!} p^l (1-p)^{(n-1)-l} \\ &= np \sum_{l=0}^{n-1} \binom{n-1}{l} p^l (1-p)^{(n-1)-l} \\ &= np, \end{aligned}$$

where in the last step we have used a standard identity discussed in more detail below (see (5.1)). △

### 5.1.1 Properties of Expected Value

Let  $X$  and  $Y$  denote any two random variables for which  $E[|X|]$  and  $E[|Y|]$  are finite. For discrete random variables with an infinite number of terms, this means the sums  $\sum_x xP(X = x)$  and  $\sum_y yP(Y = y)$  are absolutely convergent. There is an analogous meaning to this for random variables that are not discrete, and we will come to that later in the notes.

**Lemma 5.1.** *Let  $X$  and  $Y$  denote any two discrete random variables for which  $E[|X|]$  and  $E[|Y|]$  are finite. Then the following hold:*

1. For any value  $c \in \mathbb{R}$ ,  $E[cX] = cE[X]$ .
2.  $E[X + Y] = E[X] + E[Y]$ .
3. If  $X$  is a discrete random variable, and  $g$  is any real-valued function, then

$$E[g(X)] = \sum_x g(x)P(X = x),$$

*as long as the right-hand side of this expression is absolutely convergent.*

4. If  $X$  and  $Y$  are independent random variables, then  $E[XY] = E[X]E[Y]$ .

For (4), two discrete random variables  $X$  and  $Y$  are said to be independent if

$$P((X = x) \cap (Y = y)) = P(X = x)P(Y = y)$$

for all possible  $x$  and  $y$ .

**Proof of Lemma 5.1.** For (1), we define a new random variable  $Y = cX$ , and compute

$$\begin{aligned} E[Y] &= \sum_y yP(Y = y) = \sum_x cxP(cX = cx) \\ &= \sum_x cxP(X = x) = c \sum_x xP(X = x) = cE[X]. \end{aligned}$$

For (2), let's first check that for any  $X$  and  $Y$  as in the lemma, we have the relation

$$\sum_x P((X = x) \cap (Y = y)) = P(Y = y).$$

To see this, we observe that the collection of sets  $(X = x)$  necessarily forms a partition of  $S$ , so that

$$S = \bigcup_x (X = x).$$

Then

$$(Y = y) = (Y = y) \cap S = (Y = y) \cap \left( \bigcup_x (X = x) \right) = \bigcup_x ((Y = y) \cap (X = x)).$$

In this way, we see that

$$P(Y = y) = P\left(\bigcup_x ((Y = y) \cap (X = x))\right) = \sum_x P((Y = y) \cap (X = x)),$$

where we have used Axiom 3.

Now, we compute

$$\begin{aligned} E[X + Y] &= \sum_x \sum_y (x + y) P((X = x) \cap (Y = y)) \\ &= \sum_x \sum_y x P((X = x) \cap (Y = y)) + \sum_x \sum_y y P((X = x) \cap (Y = y)) \\ &= \sum_x x \sum_y P((X = x) \cap (Y = y)) + \sum_y y \sum_x P((X = x) \cap (Y = y)) \\ &= \sum_x x P(X = x) + \sum_y y P(Y = y) = E[X] + E[Y], \end{aligned}$$

where the re-ordering of summations is allowed by their absolute convergence.

For (3), we specify a new random variable  $Y = g(X)$  and compute directly from our definition of expected value,

$$E[g(X)] = \sum_{g(x)} g(x) P(g(X) = g(x)),$$

where we notice that the summation is over values that  $g(x)$  can take. We would like to sum over values that  $x$  can take, and for this let's notice that generally  $x$  can take at least as many values as  $g(x)$ . For example, if  $g(x) = x^2$  and  $x$  can take values  $-1$ ,  $0$ , and  $+1$ , then  $g(x)$  can take only values  $0$  and  $1$ . So, the sum over possible values of  $x$  will generally contain more terms than the sum over possible values of  $g(x)$ . But correspondingly the probabilities  $P(X = x)$  will be lower than the probabilities  $P(g(X) = g(x))$ . E.g., in our example with  $g(x) = x^2$ ,

$$1 \cdot P(g(X) = 1) = g(-1)P(X = -1) + g(1)P(X = +1).$$

More generally, but with the same logic, we see that

$$g(x)P(g(X) = g(x)) = \sum_{\{y:g(x)=g(y)\}} g(y)P(X = y),$$

allowing us to conclude that

$$\sum_{g(x)} g(x)P(g(X) = g(x)) = \sum_y yP(X = y) = \sum_x xP(X = x),$$

where in the final equality I'm just emphasizing that it doesn't matter what we call the variable of summation.

Finally, the proof of Item (4) is left to the homework.  $\square$

### 5.1.2 Conditional Expected Value

As with probabilities, we often want to compute an expected value, given some additional information. For example, suppose that in the two-flip game we know that the first flip lands heads. The expected value of  $X$  given this information is computed as,

$$E[X|\text{First flip heads}] = 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{2} = 1.5.$$

More generally, we have that for any event  $A$

$$E[X|A] = \sum_x xP(X = x|A).$$

Similarly as with probabilities, it's often useful to compute regular expected values by conditioning. In this regard, we have the following lemma.

**Lemma 5.2.** *Suppose the events  $\{A_k\}_{k=1}^n$  form a partition of some sample space  $S$ , and  $X$  is a random variable on  $S$  so that  $E[|X|]$  is a finite number. Then*

$$E[X] = \sum_{k=1}^n E[X|A_k]P(A_k).$$

**Proof.** We start with the right-hand side of the claim, and compute

$$\begin{aligned} \sum_{k=1}^n E[X|A_k]P(A_k) &= \sum_{k=1}^n \sum_x xP(X = x|A_k)P(A_k) \\ &= \sum_{k=1}^n \sum_x xP((X = x) \cap A_k) \\ &= \sum_x x \sum_{k=1}^n P((X = x) \cap A_k) \\ &= \sum_x xP(X = x) = E[X], \end{aligned}$$

where absolute convergence allows us to exchange the order of summation, and the second to last equality follows since

$$\bigcup_{k=1}^n (X = x) \cap A_k = (X = x) \cap S = (X = x).$$

□

**Example 5.5.** Compute the expected number of rolls of a fair pair of six-sided dice until a sum of 5 occurs.

Let

$N = \#$  of rolls until a sum of 5 occurs (a random variable)

$A =$  event a sum of 5 occurs on the first roll.

First, there are 4 ways to get a sum of 5, so  $P(A) = \frac{4}{36} = \frac{1}{9}$ . Noting also that  $A$  and  $A^c$  form a partition of  $S$ , we can use Lemma 5.2 to compute

$$\begin{aligned} E[N] &= E[N|A]P(A) + E[N|A^c]P(A^c) \\ &= 1 \cdot \frac{1}{9} + (1 + E[N]) \cdot \frac{8}{9}. \end{aligned}$$

We can now solve this algebraic equation for  $E[N]$  to see that  $E[N] = 9$ .

**Example 5.6 (The frustrated mouse).** A certain mouse is placed in the center of a maze, surrounded by three paths that open with varying widths. The first path returns him to the center after two minutes; the second path returns him to the center after four minutes; and the third path leads him out of the maze after one minute. Due to the differing widths, the mouse chooses the first path 50% of the time, the second path 30% of the time, and the third path 20% of the time. Determine the expected number of minutes it will take for the mouse to escape.

Let

$M = \#$  of minutes until the mouse escapes

$D_k =$  event the mouse takes door  $k$ .

Expected values conditioned on the events  $\{D_k\}_{k=1}^3$  are easy to calculate. For Door 3,  $E[M|D_3] = 1$ . That is, if the mouse chooses door 3, we know he will escape in 1 minute. On the other hand, for doors 1 and 2 the mouse will wander through the maze and then find himself back where he started. In particular,  $E[M|D_1] = 2 + E[M]$  and  $E[M|D_2] = 4 + E[M]$ . For example, in the case of door 1, the expected number of minutes it takes for the mouse to escape is the two minutes he spends getting back to his starting point plus his expected value of starting over. (We assume the mouse doesn't learn anything from taking the wrong doors.) Using Lemma 5.2, we find

$$\begin{aligned} E[M] &= \sum_{k=1}^3 E[M|D_k]P(D_k) \\ &= E[M|D_1] \cdot .5 + E[M|D_2] \cdot .3 + E[M|D_3] \cdot .2 \\ &= (2 + E[M]) \cdot .5 + (4 + E[M]) \cdot .3 + 1 \cdot .2, \end{aligned}$$

which is an algebraic equation that can be solved for  $E[M] = 12$ .

△

## 5.2 Variance and Covariance

Consider the following three random variables,

$$W = 0, \text{ prob } 1; \quad Y = \begin{cases} -1, & \text{prob } \frac{1}{2} \\ +1, & \text{prob } \frac{1}{2} \end{cases}; \quad Z = \begin{cases} -100, & \text{prob } \frac{1}{2} \\ +100, & \text{prob } \frac{1}{2} \end{cases}.$$

We see immediately that though these three random variables are very different, the expected value of each is the same,  $E[W] = E[Y] = E[Z] = 0$ . The problem is that the expected value of a random variable does not provide any information about how far the values the random variable takes on can deviate from one another. We measure this with *variance*, defined as

$$\text{Var}[X] = E[(X - E[X])^2].$$

That is, we study the squared difference between realizations of the random variable and the mean of the random variable. Computing directly from this definition, we find the variance of each random variable above,

$$\begin{aligned} \text{Var}[W] &= (0 - 0)^2 \cdot 1 = 0, \\ \text{Var}[Y] &= (-1 - 0)^2 \cdot \frac{1}{2} + (1 - 0)^2 \cdot \frac{1}{2} = 1, \\ \text{Var}[Z] &= (-100 - 0)^2 \cdot \frac{1}{2} + (100 - 0)^2 \cdot \frac{1}{2} = 100^2. \end{aligned}$$

Typically, a more intuitive measure of such deviation is the square root of variance, which we refer to as the *standard deviation*, and typically denote  $\sigma$ . We notice that by expanding the square in our definition of variance, we obtain the alternative expression

$$\text{Var}[X] = E[X^2] - E[X]^2.$$

**Example 5.7.** (Computing variance by conditioning) In this example, we employ a conditioning argument to determine the variance on the number of roles required in Example 5.5 for a sum of 5 to occur.

Letting  $N$  and  $A$  be as in example 5.5, we have

$$\text{Var}[N] = E[N^2] - E[N]^2,$$

for which we compute

$$\begin{aligned} E[N^2] &= E[N^2|A]P(A) + E[N^2|A^c]P(A^c) \\ &= 1 \cdot \frac{1}{9} + E[(1 + N)^2] \cdot \frac{8}{9} = \frac{1}{9} + E[1 + 2N + N^2] \frac{8}{9} \\ &= 1 + 2 \frac{8}{9} E[N] + \frac{8}{9} E[N^2]. \end{aligned}$$

We already know  $E[N]$  from Example 5.5, so we can now solve for

$$E[N^2] = 9 \cdot \left(1 + \frac{16}{9} \cdot 9\right) = 9 \cdot 17 = 153.$$

We have, finally,

$$\text{Var}[N] = 9 \cdot 17 - 9^2 = 9 \cdot 8 = 72.$$

△

**Example 5.8.** Show that the variance of a binomial random variable is  $\text{Var}[X] = np(1-p)$ .

Recall from Example 5.4 that a random variable  $X$  that takes values  $0, 1, 2, \dots, n$  is said to be a *binomial* random variable with sample size  $n$  and probability  $p$  if its probability mass function is given by

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

We showed in Example 5.4 that  $E[X] = np$ , so we need only compute  $E[X^2]$ . We have

$$\begin{aligned} E[X^2] &= \sum_{k=0}^n k^2 \binom{n}{k} p^k (1-p)^{n-k} = \sum_{k=1}^n k^2 \frac{n!}{(n-k)!k!} p^k (1-p)^{n-k} \\ &= \sum_{k=1}^{\infty} k \frac{n!}{(n-k)!(k-1)!} p^k (1-p)^{n-k}, \end{aligned}$$

where we have observed that the  $k = 0$  term gives no contribution. Setting now  $l = k - 1$  we find

$$E[X^2] = \sum_{l=0}^{n-1} (l+1)np \binom{n-1}{l} p^l (1-p)^{(n-1)-l} = npE[Y+1],$$

where  $Y$  denotes a binomial random variable that takes values  $0, 1, 2, \dots, n-1$ . I.e.,  $E[Y+1] = E[Y] + E[1] = (n-1)p + 1$ . We have, then,

$$E[X^2] = np(np - p + 1),$$

so that

$$\text{Var}[X] = E[X^2] - E[X]^2 = (np)^2 + np(1-p) - (np)^2 = np(1-p),$$

which is the claim. △

We can generalize the idea of variance to two random variables  $X$  and  $Y$ . We define the *covariance* of  $X$  and  $Y$  as

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

### 5.3 Cumulative Distribution Functions

The *cumulative distribution function*,  $F(x)$ , for a random variable  $X$  is defined for all real  $-\infty < x < +\infty$  as

$$F(x) = P(X \leq x).$$

For  $X$  as in the two-flip game above, we have

$$F(x) = \begin{cases} 0, & -\infty < x < 1 \\ 1/4, & 1 \leq x < 2 \\ 1/2, & 2 \leq x < 3 \\ 3/4, & 3 \leq x < 4 \\ 1, & 4 \leq x < \infty, \end{cases}$$

depicted graphically in Figure 5.1.

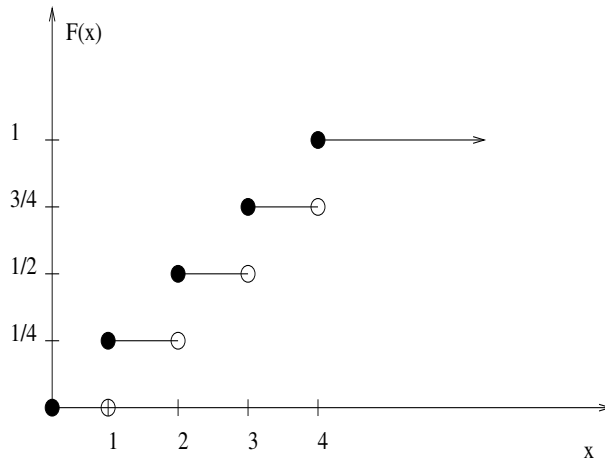


Figure 5.1: Cumulative distribution function for the two-flip game.

Below, we list for easy reference four critical properties of cumulative distribution functions,  $F(x)$ :

1.  $F(x)$  is a non-decreasing function.
2.  $\lim_{x \rightarrow +\infty} F(x) = 1$ .
3.  $\lim_{x \rightarrow -\infty} F(x) = 0$ .
4.  $F(x)$  is “right-continuous”:  $\lim_{y \rightarrow x^+} F(y) = F(x)$ .

## 5.4 Probability Mass Functions

The *probability mass function*,  $p(x)$ , for a discrete random variable  $X$  is defined by the relation

$$p(x) = P(X = x).$$

For example, in the two-flip game,  $p(1) = p(2) = p(3) = p(4) = 1/4$ . Below, we list for easy reference three critical properties of probability mass functions,  $p(x)$ .

1.  $p(x)$  is 0 except at realizations of the random variable  $X$ .
2.  $\sum_{\text{All possible } x} p(x) = 1$ .



$$3. F(y) = \sum_{x \leq y} p(x).$$

Important examples of probability mass functions include the *Poisson*, *Bernoulli*, *binomial*, and *geometric* functions.

**1. Poisson probability mass function.** A random variable  $N$  that takes values  $0, 1, 2, 3, \dots$  is said to be a *Poisson* random variable with parameter  $a$  if for some  $a > 0$  its probability mass function is given by

$$p(k) = P(N = k) = e^{-a} \frac{a^k}{k!}.$$

**2. Bernoulli probability mass function.** A random variable  $N$  that takes values  $0$  and  $1$  is said to be a *Bernoulli* random variable with probability  $p$  if its probability mass function is given by

$$p(k) = P(N = k) = \begin{cases} 1 - p, & k = 0 \\ p, & k = 1. \end{cases}$$

A single flip of a coin is a Bernoulli process with  $p = \frac{1}{2}$ .

**3. Binomial Probability mass function.** A random variable  $X$  that takes values  $0, 1, 2, \dots, n$  is said to be a *binomial* random variable with sample size  $n$  and probability  $p$  if its probability mass function is given by

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

The binomial random variable counts the number of probability  $p$  events in  $n$  trials of a Bernoulli process with probability  $p$ . For example, suppose we would like to determine the probability that 3 ones turn up in 5 rolls of a fair die. In this case,  $p = \frac{1}{6}$  (the probability of a one) and  $n = 5$ , the number of rolls. We have

$$p(3) = \binom{5}{3} \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^2 = \frac{5!}{2!3!} \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^2 = .032.$$

In order to see that the binomial probability mass function satisfies condition (2) above, we recall the binomial expansion for any integers  $a$ ,  $b$ , and  $n$ ,

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

We have, then, for the binomial probability mass function,

$$\sum_{k=0}^n p(k) = \sum_{k=0}^n \binom{n}{k} p^k (1 - p)^{n-k} = (p + (1 - p))^n = 1^n = 1. \quad (5.1)$$

**4. Geometric probability mass function.** A random variable  $X$  that takes values  $1, 2, 3, \dots$  is said to be a *geometric* random variable with probability  $p$  if its probability mass function is given by

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

A geometric random variable counts the number of trials required of a Bernoulli process with probability  $p$  to get a probability  $p$  event. In this case,  $E[X] = \frac{1}{p}$  and  $\text{Var}[X] = \frac{1-p}{p^2}$ .

## 6 Game Theory

While the modern form of game theory was initiated by the Hungarian-American mathematician John von Neumann (1903–1957) in 1928, many of the underlying ideas had been developed much earlier, dating back to the early 1700’s. A second key player in the development of game theory is the US mathematician John Nash (1928–2015), subject of the movie *A Beautiful Mind*.<sup>10</sup> Despite its possibly misleading connotations, the expression game theory has become standard terminology for the systematic study of situations of conflict and/or cooperation between players such as individuals, businesses, countries, etc. We regard a game as a situation with the following four properties<sup>11</sup>:

1. There are at least two players. A player may be an individual, but it may also be a more general entity such as a company, nation, or biological species.
2. Each player has a number of possible strategies: courses of actions that he or she may choose to follow.
3. The strategies chosen by each player determine the outcome of the game.
4. Associated to each possible outcome of the game is a collection of numerical payoffs, one to each player. These payoffs represent the value of the outcome to different players.

### 6.1 Zero-Sum Games

Suppose two players, Rose and Colin, play a game in which Rose has three possible strategies and Colin has two possible strategies. We often depict such a game in a *payoff table*, where the values in parentheses can be read as

(payoff to Rose, payoff to Colin).

**Example 6.1.** In the game depicted in Figure 6.1 if Rose plays strategy C and Colin plays strategy A then the payoff is (-5,5) so Rose loses 5 and Colin gains 5. We refer to a game such as this in which the values of each outcome sum to zero as a zero-sum game.  $\triangle$

Since Colin’s payoff is entirely determined by Rose’s in such games we often record only Rose’s payoffs. The payoff table in Figure 6.1 often appears as in Figure 6.2.

**Definition 6.1.** We refer to a two-person zero-sum game as a matrix game.

We will assume Rose and Colin are both *rational* players in the sense that they make reasonable decisions. For example, Rose would like to win 10, so she might start by playing strategy C, but as soon as Colin realizes that Rose is consistently playing strategy C he will begin playing strategy A. Of course, when Rose realizes that Colin is consistently playing strategy A she will begin playing strategy A. We can graphically depict the constantly changing choice of strategies with an arrow diagram, drawn as follows: in each column we

---

<sup>10</sup>Universal Pictures, 2001.

<sup>11</sup>Taken from [Straffin2006]

Rose/Colin	A	B
A	(2, -2)	(-3, 3)
B	(0, 0)	(2, -2)
C	(-5, 5)	(10, -10)

Figure 6.1: Payoff table for Example 6.1.

Rose/Colin	A	B
A	2	-3
B	0	2
C	-5	10

Figure 6.2: Payoff table for Example 6.1 in matrix form.

draw arrows from each entry to the largest entry in the column (this indicates what Rose's inclination will be), and in each row we draw an arrow from each entry to the smallest entry in the row (indicating what Colin's inclination will be). For example, the arrow diagram corresponding with Example 6.1 would be as depicted in Figure 6.3. In this diagram, the arrow down the second column indicates Rose's inclination for Strategy C if Colin is playing strategy B, while the arrow across the third row indicates Colin's inclination for strategy A if Rose is playing strategy C. We notice in particular that the path has no stopping point, which indicates that there is no single strategy pair that can be viewed as optimal for both Rose and Colin. We will soon return to this example, but first let's consider a few simpler games.

### 6.1.1 Dominant Strategies

The following discussion will be clarified if we have an example to refer to, so we give a convenient game in Example 6.2.

**Example 6.2.** Consider the zero-sum game depicted in Figure 6.4. In this game, Colin strategy B is clearly at least as good as Colin strategy C in each component. That is, -1 is better for Colin than 1, 1 is better than 7, 2 is better than 4, and 0 is no worse than 0. If Colin is a rational player, we can assume he will always play strategy B in favor of strategy C, and we say that his strategy B dominates his strategy C.  $\triangle$

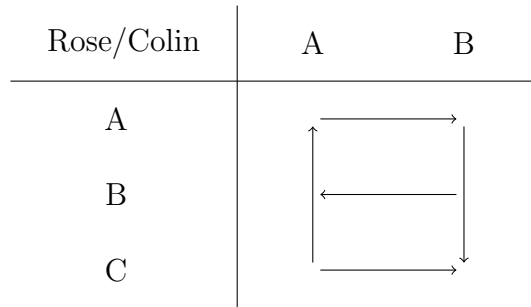


Figure 6.3: Arrow diagram for Example 6.1.

Rose/Colin	A	B	C	D
A	12	-1	1	0
B	5	1	7	-20
C	3	2	4	3
D	-16	0	0	16

Figure 6.4: Payoff table for Example 6.2.

**Definition 6.2.** A strategy  $S$  is said to dominate a strategy  $T$  if every outcome in  $S$  is at least as good as the corresponding outcome in  $T$ , and at least one outcome in  $S$  is strictly better than the corresponding outcome in  $T$ .

**Dominance Principle.** A rational player, in contest with a rational player, should never play a dominated strategy.

While this principle is straightforward, it's worth investigating just a little further. Cleverly, Colin might consider playing strategy C in the hope of inducing Rose to maximize her profits with strategy B. Then Colin will switch to strategy D and make a killing. This is where we assume Rose is also a rational player, and as such she recognizes her strategy against Colin C should be Rose C, which over time should force Colin to Colin B. According to the Dominance Principle we are justified in eliminating Colin strategy C from the table (it will never be used), and we obtain the reduced problem depicted in Figure 6.5.

We will return to Example 6.2 in the next subsection, but first we consider a game that can be solved entirely by dominance.

**Example 6.3.** Rose and Colin each put a \$1.00 ante in a pot, and then each is dealt a single card from a limitless deck consisting only of kings and aces. After looking at his card Colin must decide either to bet \$2.00 or to fold. If Colin bets then Rose must decide whether to call or fold. (We don't allow a raise in this game.) If Rose calls there is a showdown of cards

Rose/Colin	A	B	D
A	12	-1	0
B	5	1	-20
C	3	2	3
D	-16	0	16

Figure 6.5: Reduced payoff table for Example 6.2.

in which ace beats king, and in the event of a tie no money exchanges hands. Our goal is to determine optimal strategies for both Rose and Colin.

First, we want to cast this as a matrix problem. Colin's possible strategies are

Colin A: Bet only with the ace

Colin B: Always bet

Notice that Colin's strategy B allows for the possibility of bluffing in the sense that a king can't possibly be the best hand. Likewise, Rose's possible strategies are

Rose A: Call only with an ace

Rose B: Always call

We now construct a matrix game by considering Rose's expected payoff for each of the four strategy combinations.<sup>12</sup> For this game there are four possible equally likely events from the deal,

$A_1$  = event Rose gets an ace and Colin gets an ace

$A_2$  = event Rose gets an ace and Colin gets a king

$A_3$  = event Rose gets a king and Colin gets an ace

$A_4$  = event Rose gets a king and Colin gets a king

Using the fact that  $P(A_j) = \frac{1}{4}$  for all  $j = 1, 2, 3, 4$ , we can compute Rose's expected values as follows<sup>13</sup>:

*Rose A, Colin A.*

$$E[R] = 0 \cdot P(A_1) + 1 \cdot P(A_2) - 1 \cdot P(A_3) + 1 \cdot P(A_4) = \frac{1}{4}.$$

<sup>12</sup>Later, we will work with the expected value associated with each of Rose's strategies, and it's worth stressing here that that will be a fundamentally different thing.

<sup>13</sup>In these expressions we continue to write out  $P(A_j)$  to clarify the calculations.

Rose A, Colin B.

$$E[R] = 0 \cdot P(A_1) + 3 \cdot P(A_2) - 1 \cdot P(A_3) - 1 \cdot P(A_4) = \frac{1}{4}.$$

Rose B, Colin A.

$$E[R] = 0 \cdot P(A_1) + 1 \cdot P(A_2) - 3 \cdot P(A_3) + 1 \cdot P(A_4) = -\frac{1}{4}.$$

Rose B, Colin B.

$$E[R] = 0 \cdot P(A_1) + 3 \cdot P(A_2) - 3 \cdot P(A_3) + 0 \cdot P(A_4) = 0.$$

The matrix game associated with these values is given in Figure 6.6 below.

Rose/Colin	A	B
A	$\frac{1}{4}$	$\frac{1}{4}$
B	$-\frac{1}{4}$	0

Figure 6.6: Payoff table for Example 6.3.

Rose's strategy A clearly dominates her strategy B, so she should play only Strategy A, and this will ensure that she gets  $\frac{1}{4}$  per game. △

### 6.1.2 Saddle Points

Consider now the arrow diagram for the reduced game from Example 6.2, depicted in Figure 6.7. Clearly, this diagram places emphasis on the joint strategy (C,B). (In joint strategies, the row strategy will be given first, the column strategy second.) If Rose knows Colin will play B then she will certainly play C, and likewise if Colin knows Rose will play C then he will certainly play B. In this way (C,B) is regarded as an equilibrium strategy, referred to as a *saddle point* in game theory.

**Definition 6.3.** *An outcome in a matrix game (with payoffs recorded to the row player) is called a saddle point if the entry at the outcome is both its row minimum and its column maximum.*

**Definition 6.4.** *In a matrix game if there is a value  $v$  so that Rose has a strategy that guarantees that she will win at least  $v$ , and Colin has a strategy that guarantees Rose will win no more than  $v$ , then  $v$  is called the value of the game.*

**Saddle Point Principle.** *If a matrix game between two rational players has at least one saddle point, then both players should play a strategy that contains a saddle point.*

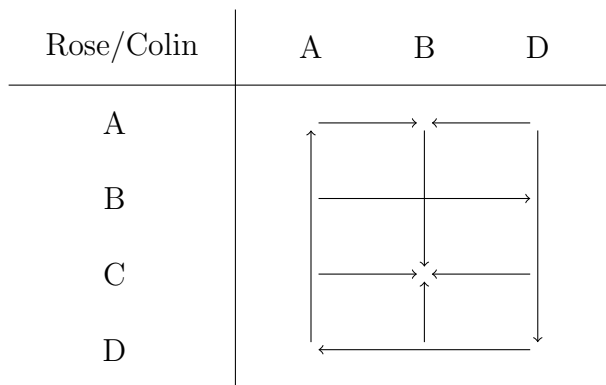


Figure 6.7: Arrow diagram for Example 6.2.

The reasoning is clear: for example, if Rose plays anything other than a saddle, Colin can ensure that she will not do as well as she would if she played a saddle. Of course a matrix game can have multiple saddle points.

**Example 6.4.** Consider the matrix game depicted in Figure 6.8. An arrow diagram for this game is shown in Figure 6.9, and it's clear from this diagram that all four ones (for joint strategies (A,A), (A,C), (D,A), and (D,C)) are saddle points..

Rose/Colin	A	B	C	D
A	1	5	1	2
B	0	-2	-5	10
C	-5	16	0	-10
D	1	13	1	2

Figure 6.8: Payoff table for Example 6.4.

**Theorem 6.1.** *We have the following:*

- (i) *Any two saddle points in a matrix game have the same value.*
- (ii) *If a matrix game has at least one saddle point its value is the value of the game.*

**Proof.** For (i) we first observe that if two saddle points are in the same row or the same column they must have identical values, because two different values can neither both be a row minimizer nor a column maximizer. Suppose, then, that  $a$  and  $d$  are any two saddle

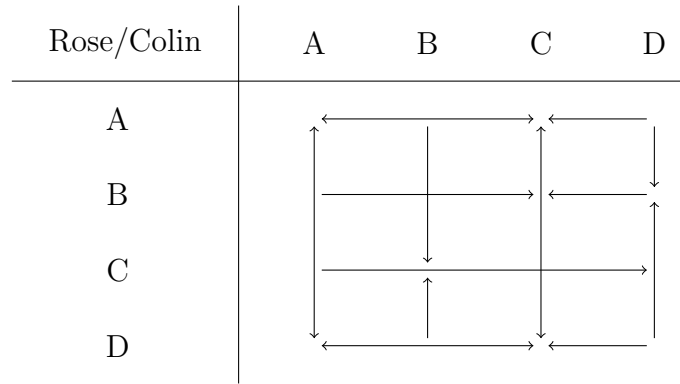


Figure 6.9: Arrow diagram for Example 6.4.

values that are not in either the same row or the same column:

$$\begin{array}{ccc}
 a & \cdots & b \\
 \vdots & \vdots & \vdots \\
 c & \cdots & d
 \end{array}$$

We know by definition of a saddle point that we have

$$\begin{aligned}
 c &\leq a \leq b \\
 b &\leq d \leq c,
 \end{aligned}$$

which gives

$$c \leq a \leq b \leq d \leq c.$$

This clearly asserts  $a = b = c = d$ .

Finally, we observe that (ii) is simply a restatement of the definition of a saddle point. Since it is the minimum value in a row Rose is guaranteed to achieve it by playing that row, and since it is the maximum value in its column Colin is guaranteed that Rose cannot do better by switching strategies.  $\square$

### 6.1.3 Minimax Method

While the graphical (arrow diagram) method for locating saddle points can work in principle on any matrix game with saddle points, it's convenient to have an analytic approach that, for example, could easily be implemented numerically for large systems. For this, we simply think through what the arrow diagram does. First, the row arrows identify the minimum value(s) across each row. Next, the column arrows identify the maximum values along each column. If an entry is *both* its row minimizer and its column maximizer it is a saddle point. This means that in locating saddle points, we can proceed as follows: (1) Find the minimum across each row and compute the maximum of these minimizers (i.e., the maximin), and (2) Find the maximum along each column and compute the minimum of these maximizers (i.e., the minimax). A value is a saddle point if and only if the maximin and the minimax are



the same value. That is, if  $v$  is a maximin then it is certainly a row minimizer, and if it is a minimax it is certainly a column maximizer. If it is both a maximin and a minimax then it is both a row minimizer and a column maximizer, and that, by definition, is a saddle point. (In fact, what this says is that any value that is in both the list of row minimizers and the list of column maximizers is a saddle, but since the smallest of the maximizers will necessarily be greater than or equal to the largest of the minimizers, the two sets of numbers can only agree if the maximin is the minimax.)

As an example, let's carry out the minimax method for the reduced version of the game in Example 6.2. In this case, we see in Figure 6.10 that the maximum of the row minima (i.e., the maximin) is 2, and likewise the minimum of the column maxima (i.e., the minimax) is 2. The corresponding saddle point is  $(C, B)$ .

Rose/Colin	A	B	D	row mins
A	12	-1	0	-1
B	5	1	-20	-20
C	3	2	3	2
D	-16	0	16	-16
column maxes	12	2	16	

Figure 6.10: Minimax method for Example 6.2.

### 6.1.4 Mixed Strategies

We have seen that in games in which there is at least one saddle point each player should play a strategy that contains a saddle point, and this process will select a particular joint strategy (or perhaps a collection of joint strategies with equivalent values for both players). In other words, we now completely understand two-person zero-sum games in which there is a saddle point. Next, we turn to the case of games in which there is no saddle point.

**Example 6.5.** Consider the matrix game depicted in Figure 6.11. In this case, it's easy to see from the arrow diagram in Figure 6.12 that there is no saddle point. In order to clarify what happens with the minimax method in the absence of a saddle point, we illustrate the procedure for this example in Figure 6.13. method below. We see that the maximin is 0 and the minimax is 2. Since the maximin does not agree with the minimax there cannot be a saddle point.

The absence of a saddle point for Example 6.5 means there is no single strategy that is optimal for both Rose and Colin, and so we must begin to think about how Rose and Colin should rationally proceed in such a situation. We begin by observing that if Colin consistently

Rose/Colin	A	B
A	2	-3
B	0	3

Figure 6.11: Payoff table for Example 6.5.

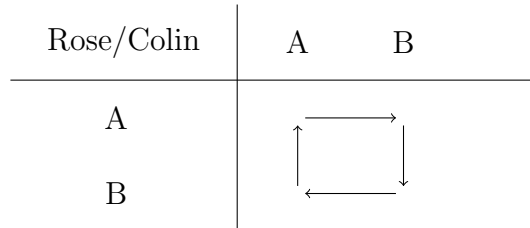


Figure 6.12: Arrow diagram for Example 6.5.

plays strategy A then Rose can profitably play strategy A, while if Colin consistently plays strategy B then Rose can profitably play strategy B. Even if Colin switches from one strategy to the other, Rose can play profitably as long as she can intelligently guess which strategy Colin will play next. In particular, notice that Rose does not have to guess correctly every time: she can play profitably as long as she can guess Colin correctly more than half the time. This suggests two things: (1) Colin should play a mixed strategy that includes both A and B, and (2) while Colin might play strategy A a different percentage of the time than he plays strategy B, he should make the choice probabilistically so that it's impossible (for all practical purposes) for Rose to guess anything about what he will play next.

In order to see how this works, let's suppose Colin plays strategy A with probability  $x$  and strategy B with probability  $1 - x$  (i.e., he always plays one or the other). In this case Rose's outcome from each of her strategies becomes a random variable, depending on what Colin plays. Let  $R_A$  denote Rose's outcome while playing strategy A and  $R_B$  denote Rose's outcome while playing strategy B. Then we have

$$E[R_A] = 2x - 3(1 - x) = 5x - 3$$

$$E[R_B] = 3(1 - x) = -3x + 3.$$

If one of these is larger than the other then Rose can profitably proceed by choosing the strategy with the larger expected outcome, so Colin should proceed by choosing  $x$  in such a way that these two expected values give the same outcome. To visualize this, we refer to the MATLAB plot in Figure 6.14 of the lines  $y = 5x - 3$  and  $y = -3x + 3$ .

In principle, Rose can always choose the strategy that gives her the expected value associated with the upper line. For example, if Colin chooses  $x = .4$ , the dashed line ( $y = -3x + 3$ ) is above the solid line, so Rose would choose strategy B. In light of this, we see that Colin's optimal choice will be the value of  $x$  where the two lines intersect. That is,

Rose/Colin	A	B	row mins
A	2	-3	-3
B	0	3	0
column maxes	2	3	

Figure 6.13: Minimax method for Example 6.5.

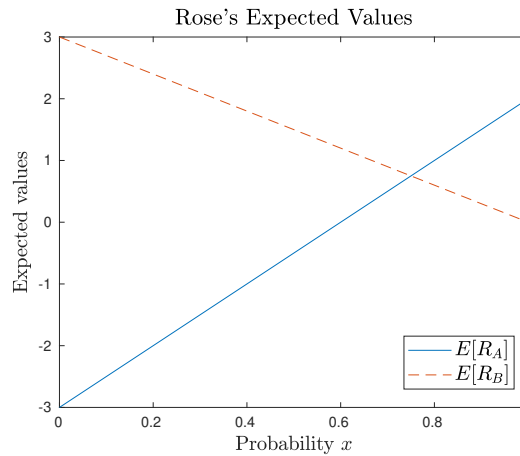


Figure 6.14: Plot of Rose's expected values for Example 6.5.

if he chooses any other value for  $x$  then Rose can certainly play a strategy that will give her a higher outcome than the intersection value. We find  $x$  by solving

$$5x - 3 = -3x + 3 \Rightarrow x = \frac{3}{4}.$$

We conclude that Colin should play strategy A  $\frac{3}{4}$  of the time and strategy B  $\frac{1}{4}$  of the time. For example, whenever it is Colin's time to make a decision he might glance at his watch. If the seconds hand is in the first quarter of the face he can play strategy B, while if it is in the latter three quarters of the face he can play strategy A. As long as Rose can't see his watch it doesn't even matter if she knows what he's doing, or even if she knows what his value of  $x$  is. If Colin plays this strategy then Rose's expectation will be the same for either of her strategies,

$$5\left(\frac{3}{4}\right) - 3 = \frac{3}{4},$$

where we emphasize that the equivalence here of Colin's probability  $x$  and Rose's expected value is coincidental.

Now of course Rose wants to think about the game in exactly the same way, so we let  $y$  denote the probability she plays strategy A and  $1 - y$  the probability she plays strategy B.

In this case Colin's expectations when playing strategies A and B are

$$\begin{aligned} E[C_A] &= -2y \\ E[C_B] &= 3y - 3(1 - y) = 6y - 3. \end{aligned}$$

(Many authors change the sign here (i.e., use the same signs as in the matrix), so that these expectations actually correspond with Rose's payoffs.) Arguing precisely as we did for Colin, we reason that Rose should equate these to expressions,

$$-2y = 6y - 3 \Rightarrow y = \frac{3}{8}.$$

We conclude that Rose should play strategy A  $\frac{3}{8}$  of the time and strategy B  $\frac{5}{8}$  of the time. Colin's outcome will be

$$-2\left(\frac{3}{8}\right) = -\frac{3}{4},$$

as we knew it should be for a zero-sum game. Notice that in our terminology  $\frac{3}{4}$  is the value of this game. △

### 6.1.5 The Method of Oddments

It's possible to analyze  $2 \times 2$  matrix games in general. Consider the general case depicted in Figure 6.15.

Rose/Colin	A	B
A	$a$	$b$
B	$c$	$d$

Figure 6.15: General  $2 \times 2$  game.

For this calculation we will assume that the game *does not* have a saddle, and we stress that this needs to be checked before the method is applied. First, suppose the two largest entries are those in the first row,  $a$  and  $b$ . Then clearly the column arrows in an arrow diagram will both point upward and the minimum of  $a$  and  $b$  will be a saddle. Likewise, if the two largest values occur in any row or column there will be a saddle. Suppose, then, that the largest values occur on diagonals, either  $a$  and  $d$  or  $b$  and  $c$ . If we let  $x$  denote the probability Colin plays strategy A so that  $1 - x$  denotes the probability he plays strategy B, Rose's general expected values become

$$\begin{aligned} E[R_A] &= ax + b(1 - x) = (a - b)x + b \\ E[R_B] &= cx + d(1 - x) = (c - d)x + d. \end{aligned}$$

Upon setting  $E[R_A] = E[R_B]$  we find

$$x = \frac{(d-b)}{(a-b) - (c-d)} = \frac{(d-b)}{(a-c) + (d-b)}.$$

Here, since we assume the largest values are diagonal from one another, we know that  $a - c$  and  $d - b$  have the same sign, and so each term in parentheses on the far right of our expression for  $x$  is positive. I.e.,

$$x = \frac{|b-d|}{|a-c| + |b-d|}.$$

This says we can proceed as follows to compute Colin's probabilities: we compute the absolute values of the differences in the columns,  $|a - c|$  and  $|b - d|$ , exchange them (after the exchange these values are referred to as Colin's *oddmments*), and write

$$P(\text{Colin plays A}) = \frac{|b-d|}{|a-c| + |b-d|}$$

$$P(\text{Colin plays B}) = \frac{|a-c|}{|a-c| + |b-d|}.$$

Proceeding similarly for Rose, we obtain a general method, described schematically for Example 6.5 in Figure 6.16.

Rose/Colin	A	B	oddmments	Rose's probs
A	2	-3	5	$\frac{3}{8}$
B	0	3	-3	$\frac{5}{8}$
oddmments	2	-6	6	2
Colin's probs	$\frac{6}{8}$	$\frac{2}{8}$		

Figure 6.16: Method of oddmments for Example 6.5.

### 6.1.6 The Minimax Theorem

Before proceeding with a few more examples, we give one of the first important results developed in the study of game theory, established by John von Neumann in 1928.

**Minimax Theorem.** Every  $m \times n$  matrix game has a solution in the following sense: there is a unique game value  $v$  and optimal strategies (pure or mixed) for Rose and Colin (i.e., for the two players in the game) so that the following hold.

(i) If Rose plays her optimal strategy her expected payoff will be greater than or equal to  $v$  no matter what Colin does.

(ii) If Colin plays his optimal strategy Rose's expected payoff will be less than or equal to  $v$  no matter what Rose does.

(iii) The solution can always be found as the solution to some square (i.e.  $k \times k$  for some positive integer  $k \leq \max(m, n)$ ) subgame of the original game.

(iv) The roles of Rose and Colin can be reversed, and the same statements obtained.

**Example 6.1 cont.** Consider again the game described in Example 6.1. The Minimax Theorem asserts that there exists a solution (game value plus optimal strategies) for this game, and that it is the solution for either a  $1 \times 1$  or  $2 \times 2$  submatrix. We have already seen that there are no saddle points for this example, and so the solution cannot be that of a  $1 \times 1$  matrix (which clearly must correspond with a saddle point). We see, then, from the Minimax Theorem that the solution will be a mixed strategy that solves one of the three possible  $2 \times 2$  matrix subgames. In order to understand which subgame we should solve, we let  $x$  denote the probability that Colin plays strategy  $A$  so that  $1 - x$  denotes the probability that he plays strategy  $B$ , and we compute Rose's expected values

$$\begin{aligned} E[R_A] &= 2x - 3(1 - x) = 5x - 3 \\ E[R_B] &= 2(1 - x) = -2x + 2 \\ E[R_C] &= -5x + 10(1 - x) = -15x + 10. \end{aligned}$$

Next, we plot these three lines together in Figure 6.17.

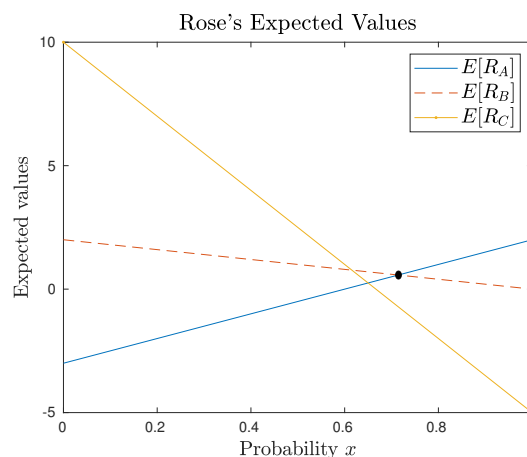


Figure 6.17: Plot for the continuation of Example 6.1.

For each value  $x$  that Colin chooses, a vertical line through  $x$  will cross all three of Rose's expected value lines. Rose can obtain the expected value associated with the largest  $y$ -value

at which these lines intersect. So Colin should choose the value  $x$  at the lowest intersection so that no lines are above the intersection. We refer to the curve created by taking the top line at each point as the upper *envelope* for these lines. In this terminology, Colin should choose  $x$  to correspond with the minimum value on the upper envelope. This point is indicated on Figure 6.17 with a dark dot. We find that the optimal strategy will be the mixed solution for Rose's strategies  $A$  and  $B$ . Graphically, the value of  $x$  has already been located, but to be precise we will solve the associated matrix game depicted in Figure 6.18.

Rose/Colin	A	B
A	2	-3
B	0	2

Figure 6.18: Payoff table for reduced game for Example 6.1.

We can solve this game by the method of oddments. Rose's oddments are respectively 2 and 5, so

$$P(\text{Rose plays A}) = \frac{2}{7}$$

$$P(\text{Rose plays B}) = \frac{5}{7},$$

and Colin's oddments are respectively 5 and 2 so

$$P(\text{Colin plays A}) = \frac{5}{7}$$

$$P(\text{Colin plays B}) = \frac{2}{7}.$$

△

**Example 6.6.** Consider the game depicted in Figure 6.19, for which we can readily see from an arrow diagram that there are no saddle points.

Rose/Colin	A	B	C	D	E
A	-2	5	1	0	-4
B	3	-3	-1	3	8

Figure 6.19: Payoff table for Example 6.6.

If we let  $y$  denote the probability that Rose plays strategy  $A$  so that  $1 - y$  denotes the probability Rose plays strategy  $B$ , then Colin's expected values are

$$\begin{aligned} E[C_A] &= 2y - 3(1 - y) = 5y - 3 \\ E[C_B] &= -5y + 3(1 - y) = -8y + 3 \\ E[C_C] &= -y + (1 - y) = -2y + 1 \\ E[C_D] &= -3(1 - y) = 3y - 3 \\ E[C_E] &= 4y - 8(1 - y) = 12y - 8. \end{aligned}$$

We depict this graphically in Figure 6.20, with values of  $y$  on the horizontal axis for convenient interpretation.

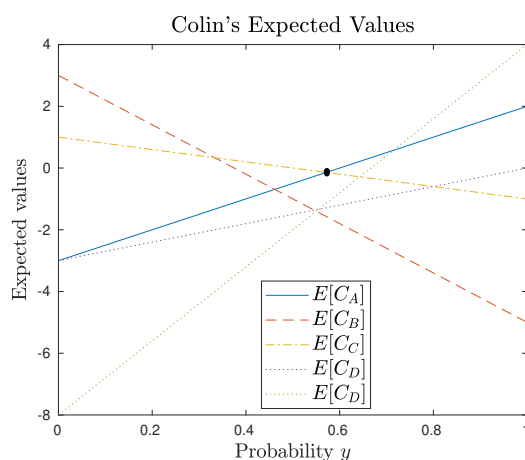


Figure 6.20: Figure for Example 6.6.

With our convention of orientation, Rose's choice of  $y$  will correspond with a vertical line through that value, and Colin will be able to bank an expected value corresponding with the largest value at which one of his expected values intersects this vertical line. In this way, Rose should choose  $y$  as the value that minimizes the upper envelope of Colin's expected value lines. This point is indicated with the black dot in Figure 6.20. We see that the optimal strategy is obtained as a solution to the  $2 \times 2$  subgame corresponding the Colin's strategies  $A$  and  $C$ . Accordingly, we must solve the game depicted in Figure 6.21.

Rose/Colin	A	C
A	-2	1
B	3	-1

Figure 6.21: Payoff table for reduced game for Example 6.6.



Using the method of oddments, we immediately find

$$P(\text{Rose plays A}) = \frac{4}{7}$$

$$P(\text{Rose plays B}) = \frac{3}{7},$$

and

$$P(\text{Colin plays A}) = \frac{2}{7}$$

$$P(\text{Colin plays C}) = \frac{5}{7}.$$

△

**Example 6.7.** Use game theory to analyze the game rock-paper-scissors.

To set this up, let's think about the game as being played for \$1.00 (paid from the loser to the winner), and let's denote the strategies  $R$ ,  $P$ , and  $S$  in the usual way. The payoff table for this game is given below, and it's straightforward to see from an arrow diagram that there are no saddle points.

Rose/Colin	R	P	S
R	0	-1	1
P	1	0	-1
S	-1	1	0

Figure 6.22: Payoff table for rock-paper-scissors.

In order to analyze the game, let's let  $x$  denote the probability that Colin plays  $R$ ,  $y$  the probability that he plays  $P$ , leaving  $1 - x - y$  as the probability he plays  $S$ . If  $R_R$ ,  $R_P$ , and  $R_S$  respectively denote Rose's payoffs while playing strategies  $R$ ,  $P$ , and  $S$ , then we can write

$$E[R_R] = 0 \cdot x - y + (1 - x - y) = -x - 2y + 1$$

$$E[R_P] = x + 0 \cdot y - (1 - x - y) = 2x + y - 1$$

$$E[R_S] = -x + y.$$

This is a bit harder to graph, but there are three planes, and Colin can minimize Rose's possible expectation by finding the point where they all intersect. I.e., we require

$$E[R_R] = E[R_P] = E[R_S],$$

and we can view this as a system of two equations for the two unknowns  $x$  and  $y$ . If we take  $E[R_R] = E[R_P]$  and  $E[R_P] = E[R_S]$ , we have respectively

$$\begin{aligned} -x - 2y + 1 &= 2x + y - 1 \\ 2x + y - 1 &= -x + y. \end{aligned}$$

Rearranging, we see that

$$\begin{aligned} 3x + 3y &= 2 \\ 3x &= 1, \end{aligned}$$

from which we conclude that  $x = \frac{1}{3}$  and  $y = \frac{1}{3}$ . As expected, we conclude that Colin should should play each strategy with equal probability, and by symmetry Rose should do the same.  $\triangle$

**Example 6.8.** Find Rose's optimal strategy for the matrix game depicted in Figure 6.23.

Rose/Colin	A	B	C
A	7	-5	1
B	-3	0	4
C	-1	6	-2

Figure 6.23: Payoff table for Example 6.8.

In order to find Rose's optimal strategy, we let  $x$  denote the probability she plays strategy  $A$ , and  $y$  the probability she plays strategy  $B$ , leaving  $1 - x - y$  to be the probability that she plays strategy  $C$ . If  $C_A$ ,  $C_B$ , and  $C_C$  respectively denote Colin's winnings playing strategies  $A$ ,  $B$ , and  $C$ , then we can write

$$\begin{aligned} E[C_A] &= -7x + 3y + (1 - x - y) = -8x + 2y + 1 \\ E[C_B] &= 5x + 0 \cdot y - 6(1 - x - y) = 11x + 6y - 6 \\ E[C_C] &= -x - 4y + 2(1 - x - y) = -3x - 6y + 2. \end{aligned}$$

Setting  $E[C_A] = E[C_B]$  and  $E[C_B] = E[C_C]$ , we see that

$$\begin{aligned} -8x + 2y + 1 &= 11x + 6y - 6 \\ 11x + 6y - 6 &= -3x - 6y + 2, \end{aligned}$$

which we can rearrange as

$$\begin{aligned} 19x + 4y &= 7 \\ 14x + 12y &= 8. \end{aligned}$$

If we multiply the first equation by 3 and subtract the second equation, we get to

$$43x = 13 \implies x = \frac{13}{43}.$$

It follows that  $y = \frac{27}{86}$ , so Rose should play strategy  $A$  with probability  $\frac{13}{43}$ , strategy  $B$  with probability  $\frac{27}{86}$ , and strategy  $C$  with probability  $\frac{33}{86}$ .  $\triangle$

## 7 Continuous Random Variables

For some random variables, the collection of possible realizations can be regarded as continuous. For example, the time between scores in a soccer match or the price of a stock at a given time are such random variables.

### 7.1 Cumulative Distribution Functions

Similarly as with discrete random variables, the cumulative distribution function,  $F(x)$ , for a continuous random variable  $X$  is defined by  $F(x) = P(X \leq x)$ .

**Example 7.1.** Suppose  $U$  is a random variable that takes values on the interval  $[0, 1]$ , and suppose additionally that we have no reason to believe that the probability of  $U$  taking any one value is different from the probability that it will take any other. We refer to such a random variable as *uniformly* distributed on  $[0, 1]$ . Let's write down an expression for the cumulative distribution function of  $U$ .

Since  $U$  is equally likely to take on any value in the interval  $[0, 1]$ , we observe that it has the same likelihood of being above  $1/2$  as being below. That is,

$$F\left(\frac{1}{2}\right) = P(U \leq 1/2) = 1/2.$$

Similarly,

$$F\left(\frac{1}{3}\right) = P(U \leq 1/3) = 1/3$$

and so on,

$$F\left(\frac{1}{n}\right) = P(U \leq 1/n) = 1/n. \tag{7.1}$$

In general, we have the relationship that if  $F$  is the CDF associated with  $U$ , then

$$F(x) = \begin{cases} 0, & x \leq 0 \\ x, & 0 \leq x \leq 1 \\ 1, & x \geq 1. \end{cases}$$

Finally, we observe that if we take the limit in (7.1) as  $n \rightarrow \infty$ , we get

$$P(U = 0) = 0.$$



- $f(x) \geq 0$  for all  $x \in \mathbb{R}$ ;
- $\int_{-\infty}^{+\infty} f(x)dx = 1$ ;
- For any interval  $[a, b] \subset \mathbb{R}$ ,  $a < b$ ,

$$P(a < X < b) = \int_a^b f(x)dx,$$

where it's assumed in this last item that the integral makes sense for all such intervals.

**Remark 7.1.** This is not the most general definition of either a continuous random variable or a probability density function, but it's all we will need for these notes. For the integral condition, it's fine to think of  $f$  as being Riemann integrable on all intervals  $[a, b]$ , and even better to think of it as being Lebesgue measurable on all such intervals.

### 7.3 Expected Value, Variance, and Covariance

**Definition 7.2.** For a continuous random variable  $X$  with probability density function  $f$  as specified in Definition 7.1, we say that  $X$  is integrable provided

$$\int_{-\infty}^{+\infty} |x|f(x)dx < \infty.$$

For integrable random variables, we define the expected value  $E[X]$  to be

$$E[X] := \int_{-\infty}^{+\infty} xf(x)dx.$$

Likewise, we say that a piecewise continuous function  $g$  is integrable with respect to  $f$  provided

$$\int_{-\infty}^{+\infty} |g(x)|f(x)dx < \infty,$$

and in this case we define the expected value of  $g(X)$  to be

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

**Definition 7.3.** If  $X$  and  $X^2$  are both integrable random variables with respect to the PDF  $f$  for  $X$ , then we define the variance of  $X$  to be

$$\text{Var}[X] := E[(X - E[X])^2] = \int_{-\infty}^{+\infty} (x - E[X])^2 dx.$$

## 7.4 Identifying Probability Density Functions

A critical issue in the study of probability and statistics regards determining the probability density function for a given random variable. We typically begin this determination through consideration of a histogram, simply a bar graph indicating the number of realizations that fall into each of a predetermined set of intervals.

**Example 7.2.** Suppose we are hired by Lights, Inc. to study the lifetime of lightbulbs (a continuous random variable). We watch 100 bulbs and record times to failure, organizing them into the following convenient ranges (in hours):

Time range	0–400	400–500	500–600	600–700	700–800	800–900
# Failed	0	2	3	5	10	10
900–1000	1000–1100	1100–1200	1200–1300	1300–1400	1400–1500	1500–1600
20	20	10	10	5	3	2

Table 7.1: Data for Lights, Inc. example.

This data is recorded in the MATLAB M-file *lights.m*.

```
%LIGHTS: Script file that defines times T at which
%lightbulbs failed.
T=[401 402 501 502 503 601 602 603 604 605 701 702 703 704 705 706 ...
707 708 709 710 801 802 803 804 805 806 807 808 809 810 901 902 903 ...
904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 ...
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 ...
1015 1016 1017 1018 1019 1020 1101 1102 1103 1104 1105 1106 1107 1108 ...
1109 1110 1201 1202 1203 1204 1205 1206 1207 1208 1209 1210 1301 1302 ...
1303 1304 1305 1401 1402 1403 1501 1502];
```

Of course, we could analyze the times more carefully in intervals of 50 hours or 10 hours etc., but for the purposes of this example, intervals of 100 hours will suffice. We define now the function

$$f(x) = \begin{cases} 0, & 0 \leq x \leq 400 \\ .0002, & 400 \leq x \leq 500 \\ .0003, & 500 \leq x \leq 600 \\ .0005, & 600 \leq x \leq 700 \\ \vdots & \vdots \\ .0002, & 1500 \leq x \leq 1600 \end{cases},$$

Let  $T$  denote the time to failure of the lightbulbs. Then we can compute the probability that a lightbulb will fail in some interval  $[a, b]$  by integrating  $f(x)$  over that interval. For example,

$$P(400 \leq T \leq 500) = \int_{400}^{500} f(x) dx = \int_{400}^{500} .0002 dx = .02.$$

$$P(600 \leq T \leq 800) = \int_{600}^{800} f(x) dx = \int_{600}^{700} .0005 dx + \int_{700}^{800} .001 dx = .15.$$

Recalling our properties of the probability density function, we see that  $f(x)$  is an approximation to the PDF for the random variable  $T$ . (It's not precise, because it is only precisely accurate on these particular intervals, and a PDF should be accurate on all intervals.)

The function  $f(x)$  is a histogram for our data, scaled by a value of 10,000 to turn numerical counts into probabilities (more on this scale in a minute). If the vector  $T$  contains all the failure times for these lights, then the MATLAB command `hist(T,12)` creates a histogram with twelve *bins* (bars) (see Figure 7.1). Notice that this histogram is precisely a scaled version of

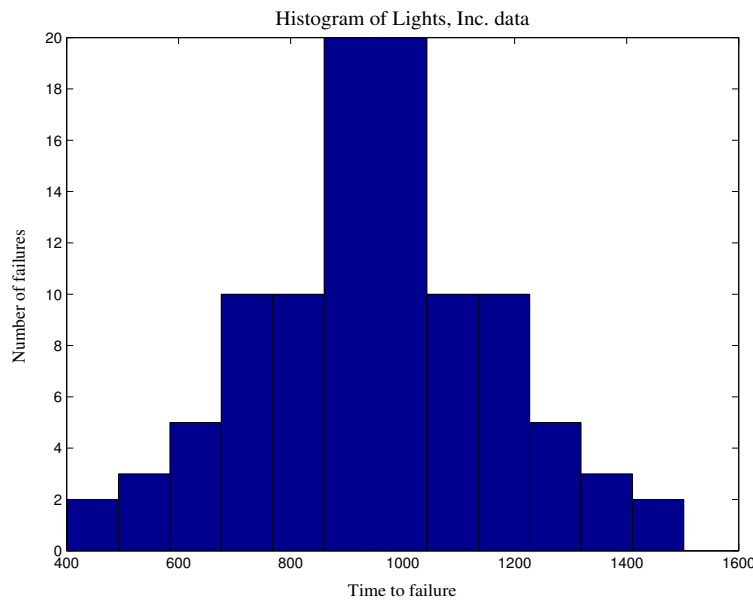


Figure 7.1: Histogram of data from Example 7.2.

$f(x)$ . As for the scaling, we choose it so that  $\int_{-\infty}^{+\infty} f(x)dx = 1$ , which can be accomplished by dividing the height of each bar in the histogram by  $\text{binwidth} \times \text{Total number of data points}$ . Here, we divide by  $100 \times 100 = 10000$ , giving  $f(x)$ .  $\triangle$

## 7.5 Useful Probability Density Functions

We typically proceed by looking at a histogram of our data, and trying to match it to the form of a smooth probability density function, and preferably one that is easy to work with. (Though as computing power becomes better and better, researchers are becoming less concerned with ease of computation.) In these notes, we will focus on the following distributions.

1. Gaussian
2. Uniform
3. Exponential
4. Weibull

- 5. Beta
- 6. Gamma
- 7. Mixture

**1. Gaussian distribution.** One of the most common probability density functions is the *Gaussian* distribution—also known as the normal distribution. The Gaussian probability density function for a random variable  $X$  with mean  $\mu$  ( $E[X] = \mu$ ) and standard deviation  $\sigma$  ( $\sigma^2 = \text{Var}[X]$ ), takes the form

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

**Example 7.2 cont.** For the case of our Lights, Inc. data we compute  $f(x)$  with the following MATLAB script.

```
>>mu=mean(T)
mu =
956.8800
>>sd=std(T)
sd =
234.6864
>>x=linspace(400,1600,100);
>>f=1/(sqrt(2*pi)*sd)*exp(-(x-mu).^2/(2*sd^2));
>>plot(x,f,'-')
```

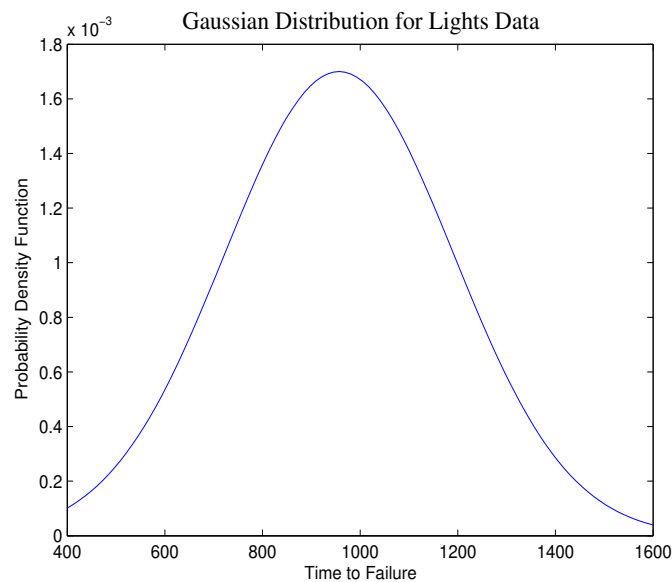


Figure 7.2: Gaussian distribution for Lights, Inc. data.



In order to compare our fit with our data, we will have to scale the Gaussian distribution so that it is roughly the same size as the histogram. More precisely, we know that the Gaussian distribution must integrate to 1 while an integral over our histogram is given by a sum of the areas of its rectangles. In the event that these rectangles all have the same width, which is the generic case for histograms, this area is

$$\text{Total area} = [\text{width of a single bar}] \times [\text{total points}].$$

In MATLAB, if the histogram command  $[n,c]=\text{hist}(T,12)$  will return a vector  $n$  containing the number of data points in each bin, and another vector  $c$  containing the center point of each bin. The binwidth can be computed from  $c$  as, for example,  $c(2) - c(1)$ . The scaling can be computed, then, as

$$\text{Scale} = (c(2) - c(1)) * \text{sum}(n).$$

Assuming  $x$ ,  $\mu$  and  $sd$  are defined as above, we use the following MATLAB script.

```
>>[n,c]=hist(T);
>>hist(T)
>>f=(c(2)-c(1))*sum(n)/(sqrt(2*pi)*sd)*exp(-(x-mu).^2/(2*sd^2));
>>hold on
>>plot(x,f,'r')
```

(see Figure 7.3).

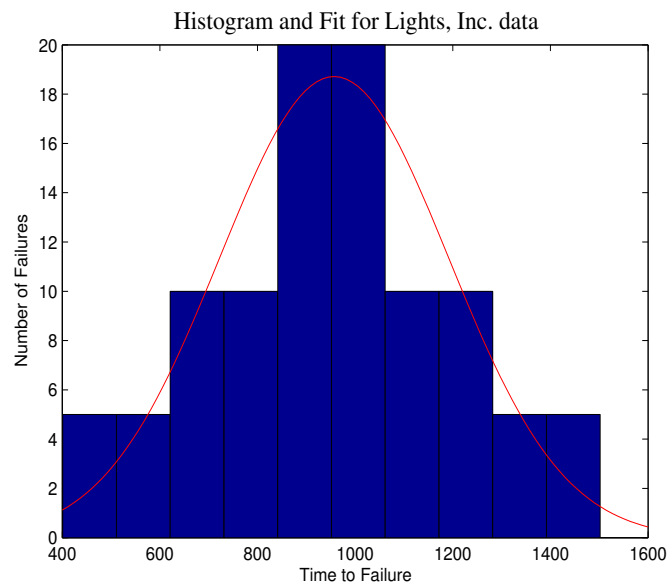


Figure 7.3: Lights, Inc. data with Gaussian distribution.

The Gaussian distribution is typically useful when the values a random variable takes are clustered near its mean, with the probability that the value falls below the mean equivalent to the probability that it falls above the mean. Typical examples include the height of a

randomly selected man or woman, the grade of a randomly selected student, and the velocity of a molecule of gas.

**2. Uniform Distribution.** The *uniform* probability density function has the form

$$f(x) = \begin{cases} \frac{1}{b-a}, & a < x < b \\ 0, & \text{otherwise} \end{cases},$$

and is applicable to situations for which all outcomes on some interval  $[a, b]$  are equally likely. The mean of a uniformly distributed random variable is  $\frac{a+b}{2}$ , while the variance is  $\frac{(b-a)^2}{12}$ .

**Example 7.3.** Consider the game of American roulette, in which a large wheel with 38 slots is spun in one direction and a small white ball is spun in a groove at the top of the wheel in the opposite direction. Though Newtonian mechanics could ideally describe the outcome of roulette exactly, the final groove in which the ball lands is for all practical purposes random. Of course, roulette is a discrete process, but its probability density function can be well approximated by the uniform distribution. First, we create a vector  $R$  that contains the outcomes of 5000 spins:

```
>>R=ceil(rand([5000,1])*38);
```

The following MATLAB code compares a histogram of this data with its probability density function (see Figure 7.4).

```
>>[n,c]=hist(R,38)
>>hist(R,38)
>>x=linspace(1,39,25);
>>f=(c(2)-c(1))*sum(n)*sign(x);
>>hold on
>>plot(x,f,'-')
```

**3. Exponential distribution.** The *exponential* probability density function is given by

$$f(x) = \begin{cases} ae^{-ax}, & x > 0 \\ 0, & x < 0 \end{cases},$$

where  $a > 0$ . This distribution is often employed as a model for a random variable associated with time in situations in which the time remaining until the next event is independent of the time since the previous event. Examples include the time between goals in a soccer match and the time between arrivals in a waiting line. The mean of an exponentially distributed random variable is  $1/a$ , while the variance is  $1/a^2$ .

**Example 7.4.** Consider the number of rolls between sixes on a fair die, where two sixes in a row correspond with zero roles. The M-file *roles1.m* creates a vector  $R$  containing 10,000 realizations of this random variable.

```
%ROLES1: Creates a list R of number of roles of
%a six-sided die between occurrences of a 6.
```

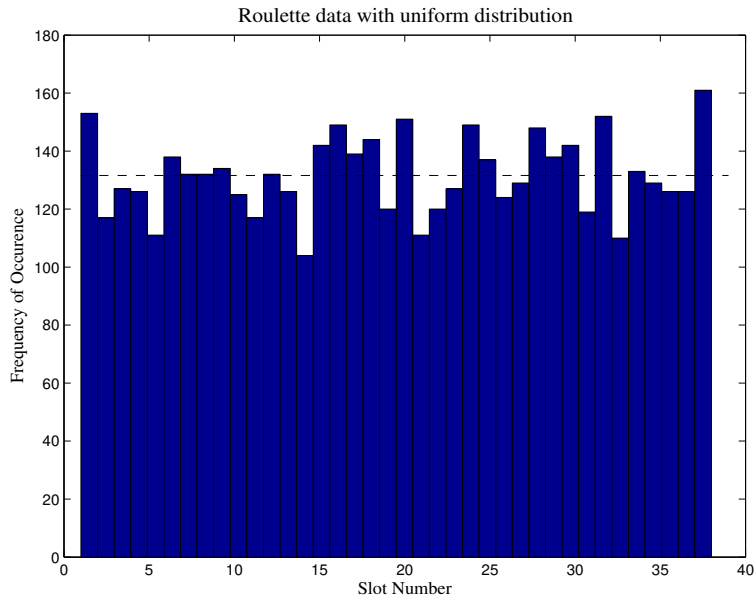


Figure 7.4: Uniform distribution with roulette data.

```

N=10000; %Number of sixes
clear R;
for k=1:N
m = 0; %Number of roles since last 6 (0 for 2 sixes in a row)
num = rand*6; %Random number between 1 and 6.
while num <= 5
m = m + 1;
num = rand*6; %Next role
end
R(k) = m;
end

```

The following MATLAB code produces Figure 7.5.

```

>>[n,c]=hist(R,max(R)+1)
>>hist(R,max(R)+1)
>>mu=mean(R)
mu =
4.9493
>>x=linspace(0,max(R),max(R));
>>f=(c(2)-c(1))*sum(n)*(1/mu)*exp(-x/mu);
>>hold on
>>plot(x,f,'-')

```

**4. Weibull Distribution.** The probability density function for the *Weibull* distribution is

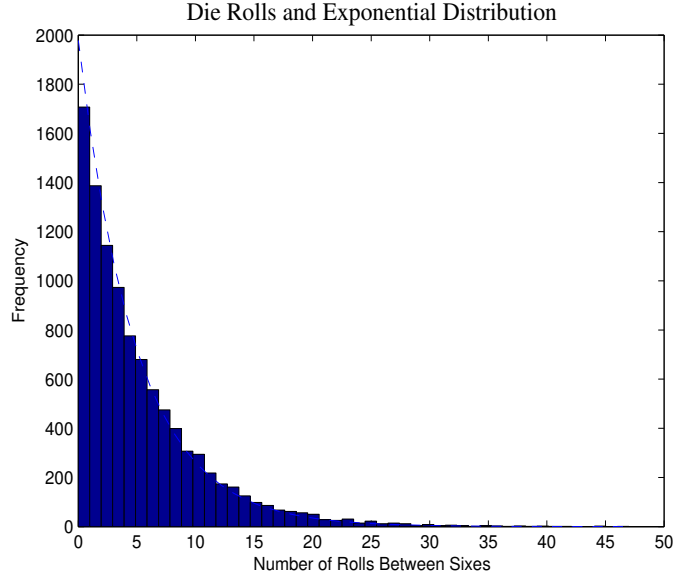


Figure 7.5: Histogram and exponential pdf for Example 7.4.

given by

$$f(x) = \begin{cases} \lambda^\beta \beta x^{\beta-1} e^{-(\lambda x)^\beta} & x > 0 \\ 0 & x < 0, \end{cases}$$

where  $\lambda > 0$  and  $\beta > 0$ , with mean and variance

$$E[X] = \frac{1}{\lambda} \Gamma\left(1 + \frac{1}{\beta}\right), \quad \text{and} \quad \text{Var}[X] = \frac{1}{\lambda^2} \left( \Gamma\left(1 + \frac{2}{\beta}\right) - \Gamma\left(1 + \frac{1}{\beta}\right)^2 \right).$$

Here,  $\Gamma(\cdot)$  denotes the *gamma* function,

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt.$$

(MATLAB has a built-in gamma function, *gamma()*.) Observe that in the case  $\beta = 1$  the Weibull distribution reduces to the exponential distribution. Named for the Swedish mechanical engineer Waloddi Weibull (1887–1979) who first suggested it, the Weibull distribution is widely used as a model for times to failure; for example, in the case of automotive parts. The Weibull probability density function for  $\beta = 2$ ,  $\lambda = 1$  and for  $\beta = \frac{1}{2}$ ,  $\lambda = 1$  is depicted in Figure 7.6.

**5. Beta Distribution.** The probability density function for the *beta* distribution is given by

$$f(x) = \begin{cases} \frac{1}{\beta(a,b)} x^{a-1} (1-x)^{b-1} & 0 < x < 1 \\ 0 & \text{otherwise,} \end{cases}$$

where  $a > 0$  and  $b > 0$ , and where the *beta* function is defined as

$$\beta(a,b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx.$$

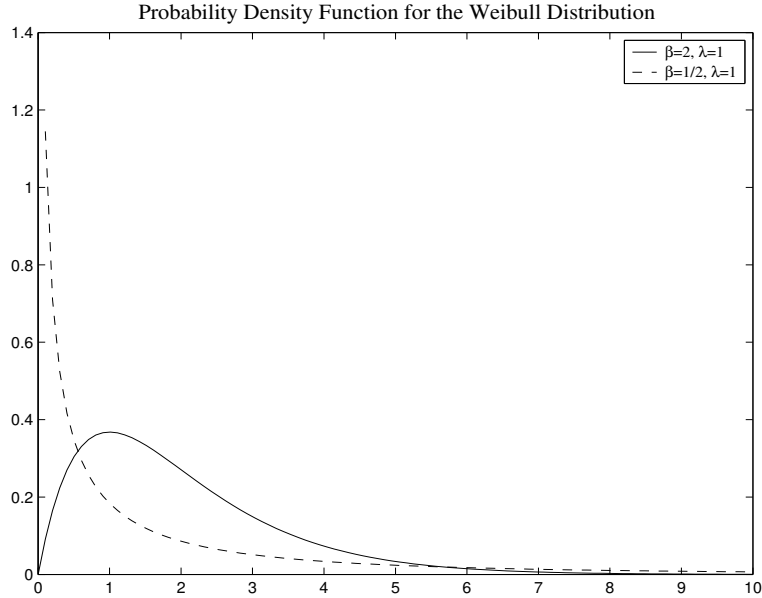


Figure 7.6: Probability density function for the Weibull distribution.

(MATLAB has a built-in beta function,  $\text{beta}(\cdot)$ .) The expected value and variance for beta random variables,  $X$ , are

$$E[X] = \frac{a}{a+b}; \quad \text{Var}[X] = \frac{ab}{(a+b)^2(a+b+1)}.$$

The beta random variable is useful in the event of slow tails; that is, when the probability density function decays at algebraic rate rather than exponential. The beta distribution for values  $a = 2$ ,  $b = 4$  and for  $a = \frac{1}{2}$ ,  $b = 2$  are depicted in Figure 7.7.

**6. Gamma Distribution.** The probability density function for the gamma distribution is given by

$$f(x) = \begin{cases} \frac{\lambda e^{-\lambda x} (\lambda x)^{n-1}}{\Gamma(n)}, & x \geq 0 \\ 0, & x < 0, \end{cases}$$

for some  $\lambda > 0$  and  $n > 0$ , where  $\Gamma(\cdot)$  denotes the gamma function as defined in the discussion of the Weibull distribution above. The mean and variance of the gamma distribution are

$$E[X] = \frac{n}{\lambda}; \quad \text{Var}[X] = \frac{n}{\lambda^2}.$$

When  $n$  is an integer, the gamma distribution is the distribution of the sum of  $n$  independent exponential random variables with parameter  $\lambda$ . The case  $\lambda = 1$ ,  $n = 2$  is depicted in Figure 7.8.

**7. Mixture Distributions.** Often, a random phenomenon will be divided into two or more characteristic behaviors. For example, if the random variable  $T$  represents service time in a certain coffee house, the time for specialty drinks may satisfy an entirely different distribution than the time for drip coffee.

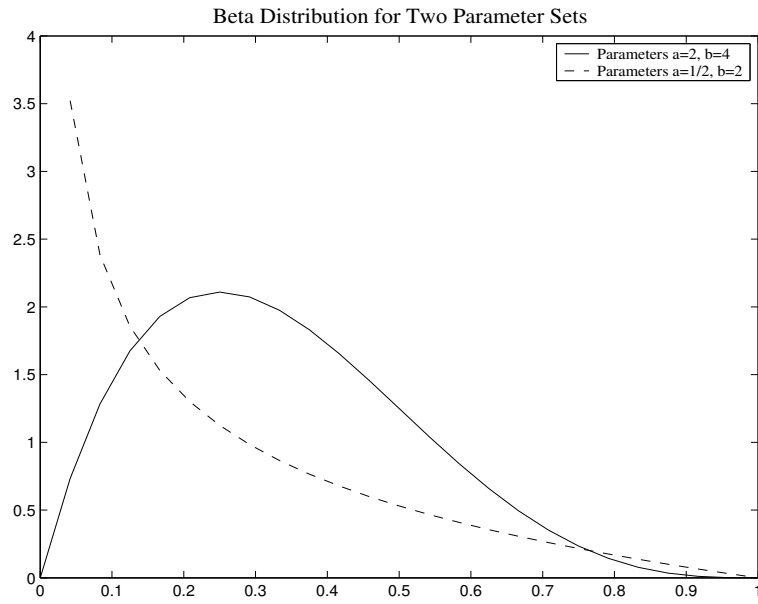


Figure 7.7: Probability Density Functions for Beta distribution.

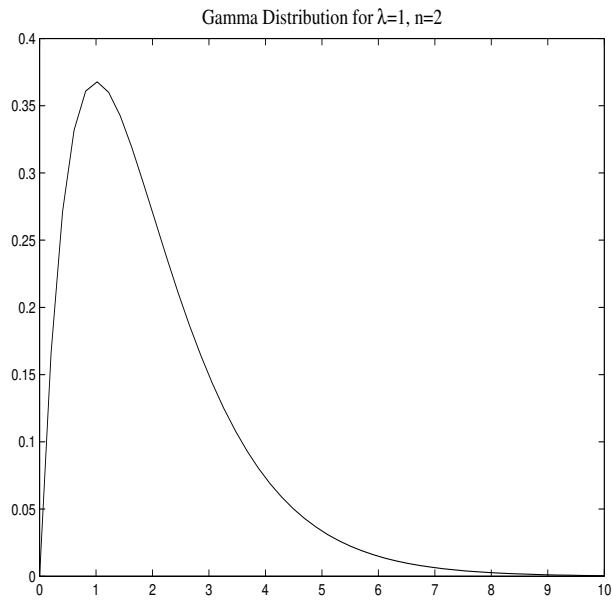


Figure 7.8: Probability density function for the gamma distribution.

**Example 7.5.** Consider a game that begins with two flips of a fair coin. If the two flips are both heads, then a coin is flipped 50 times and \$1.00 is paid for each tail, \$2.00 for each head. On the other hand, if the initial two flips are not both heads, then a fair six-sided die is rolled fifty times and \$1.00 is paid for each 1, \$2.00 is paid for each 2, \$3.00 for each 3 etc.

In the MATLAB M-file *mixture.m* just below, we compute a vector  $R$  containing as its elements the outcomes of 1,000 simulated plays of the game. (For a discussion of simulating random variables, see Section 9.)

```
%MIXTURE: Script file to run example of a mixture
%random variable.
%Experiment: Flip a coin twice. If HH, flip the coin
%fifty more times and make $1.00 for each head, and
%$2.00 for each tail. If not HH, role a fair
%six-sided die 50 times and make $1.00 for
%each 1, $2.00 for each 2, etc.
N = 1000; %Number of times to run experiment.
for j=1:N
m=0;
if rand <= .25 %Flip two heads
for k=1:50
m = m + round(rand*2+.5);
end
else
for k=1:50
m = m + round(rand*6+.5);
end
end
R(j) = m;
end
```

The MATLAB command  $hist(R, max(R))$  creates the histogram in Figure 7.9.

As expected, the payoff from the fifty coin flips satisfies a distribution entirely different from that of the payoff from the fifty die rolls. In order to analyze this data, we first need to split the data into two parts, one associated with the coin flips and the other associated with the die rolls. Generally, the two (or more) distributions will run into one another, so this step can be quite difficult, but here it is clear that we should take one set of data for payoffs below, say, 110, and the other set of data above 110. We will refer to the former as  $S$  for *small* and the latter as  $B$  for *big*. Letting  $\mu_s$  and  $\sigma_s$  represent the mean and standard deviation for  $S$ , and letting  $\mu_b$  and  $\sigma_b$  represent the mean and standard deviation for  $B$ , we will fit each clump of data separately to a Gaussian distribution. We have,

$$f_s(x) = \frac{1}{\sqrt{2\pi}\sigma_s} e^{-\frac{(x-\mu_s)^2}{2\sigma_s^2}}; \quad f_b(x) = \frac{1}{\sqrt{2\pi}\sigma_b} e^{-\frac{(x-\mu_b)^2}{2\sigma_b^2}}.$$

In combining these into a single mixture distribution, we must be sure the integral of the final distribution over  $\mathbb{R}$  is 1. (An integral over  $f_s + f_b$  is clearly 2.) Letting  $p$  represent the

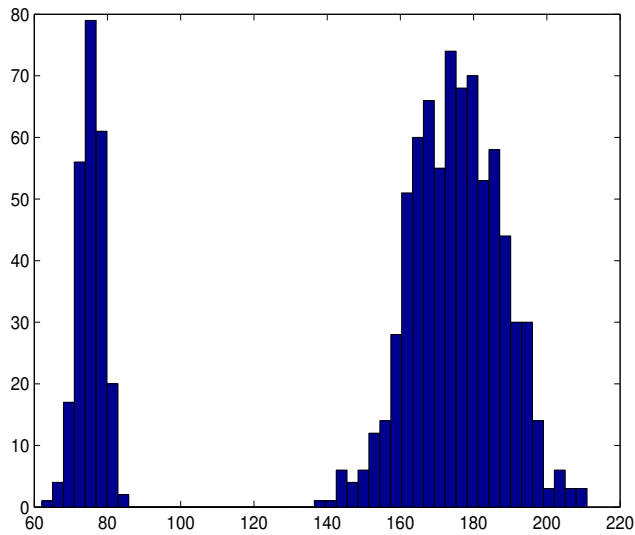


Figure 7.9: Histogram with data from Example 7.5.

probability that a play of our game falls into  $S$ , we define our mixture distribution by

$$f(x) = pf_s(x) + (1 - p)f_b(x),$$

for which

$$\begin{aligned} \int_{-\infty}^{+\infty} f(x)dx &= \int_{-\infty}^{+\infty} (pf_s(x) + (1 - p)f_b(x))dx \\ &= p \int_{-\infty}^{+\infty} f_s(x)dx + (1 - p) \int_{-\infty}^{+\infty} f_b(x) = p + (1 - p) = 1. \end{aligned}$$

We first parse our data into two sets,  $S$  for the smaller numbers and  $B$  for the bigger numbers.

```
%MIX1: Companion file for mixture.m, cleans up the data.
global N;
global R;
i = 0;
l = 0;
for k=1:N
if R(k) <= 110
i = i + 1;
S(i) = R(k);
else
l = l + 1;
B(l) = R(k);
end
end
```



The following M-file now creates Figure 7.10.

```
%MIX1PLOT: MATLAB script M-file for comparing mixture
%distribution with histogram for data created in mixture.m
hist(R,50);
mus = mean(S); sds = std(S); mub = mean(B); sdb = std(B);
p = length(S)/(length(S)+length(B));
x = linspace(0, max(B), max(B));
fs = 1/(sqrt(2*pi)*sds)*exp(-(x-mus).^2/(2*sds^2));
fb = 1/(sqrt(2*pi)*sdb)*exp(-(x-mub).^2/(2*sdb^2));
[n,c]=hist(R, 50);
f = sum(n)*(c(2)-c(1))*(p*fs+(1-p)*fb);
hold on
plot(x,f,'r')
```

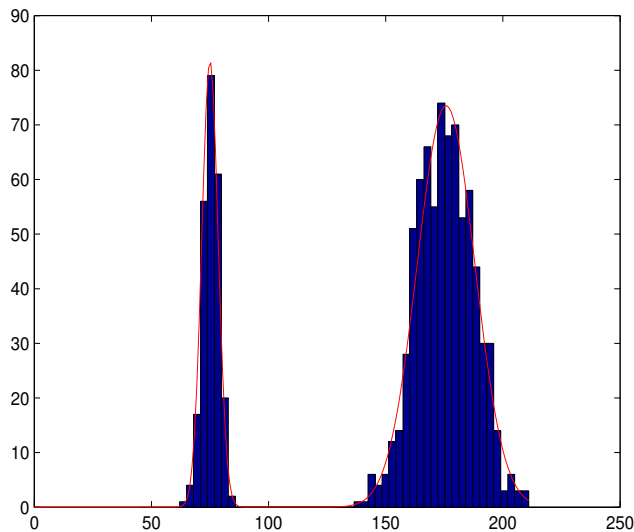


Figure 7.10: Mixture distribution with data from Example 7.5.

## 7.6 More Probability Density Functions

In this section, we list for convenient reference ten additional PDFs, though detailed discussions are omitted. We continue our numbering scheme from Section 7.5.

**8. Cauchy Distribution.** The PDF for the Cauchy distribution is

$$f(x) = \frac{1}{\pi\beta\left[1 + \left(\frac{x-\alpha}{\beta}\right)^2\right]},$$

where  $-\infty < \alpha < \infty$ , and  $\beta > 0$ . The expected value and variance for the Cauchy distribution are infinite.

**9. Lognormal Distribution.** The PDF for the lognormal distribution is

$$f(x) = \frac{1}{x\sqrt{2\pi\sigma^2}} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}}, \quad x > 0,$$

where  $-\infty < \mu < \infty$  and  $\sigma > 0$ . The lognormal distribution arises through exponentiation of the Gaussian distribution. That is, if  $G$  is a Gaussian random variable with mean  $\mu$  and variance  $\sigma^2$ , then the random variable  $X = e^G$  has the PDF given above. In this case, we have

$$\begin{aligned} E[X] &= e^{\mu + \frac{1}{2}\sigma^2} \\ \text{Var}[X] &= e^{2\mu + 2\sigma^2} - e^{2\mu + \sigma^2}. \end{aligned}$$

This random variable plays a fundamental role in the modeling of stock prices.

**10. Double Exponential (or Laplace) Distribution.** The PDF for the double exponential distribution is

$$f(x) = \frac{1}{2\beta} e^{-\frac{|x-\alpha|}{\beta}},$$

where  $-\infty < \alpha < \infty$  and  $\beta > 0$ . If  $X$  has a double exponential distribution, then

$$\begin{aligned} E[X] &= \alpha \\ \text{Var}[X] &= 2\beta^2. \end{aligned}$$

**11. Logistic Distribution.** The PDF for the logistic distribution is

$$F(x) = \frac{1}{\beta} \frac{e^{-\frac{x-\alpha}{\beta}}}{(1 + e^{-\frac{x-\alpha}{\beta}})^2},$$

where  $-\infty < \alpha < \infty$  and  $\beta > 0$ . If  $X$  has a logistic distribution, then

$$\begin{aligned} E[X] &= \alpha \\ \text{Var}[X] &= \frac{\beta^2\pi^2}{3}. \end{aligned}$$

**12. Rayleigh Distribution.** The PDF for the Rayleigh distribution is

$$f(x) = \frac{x}{\sigma^2} e^{-\frac{x^2}{2\sigma^2}}, \quad x > 0$$

where  $\sigma > 0$ . If  $X$  has a Rayleigh distribution, then

$$\begin{aligned} E[X] &= \sigma\sqrt{\frac{\pi}{2}} \\ \text{Var}[X] &= 2\sigma^2\left(1 - \frac{\pi}{4}\right). \end{aligned}$$

The Rayleigh distribution is used in the modeling of communications systems and in reliability theory.

**13. Pareto Distribution.** The PDF for the Pareto distribution is

$$f(x) = \frac{\beta\alpha^\beta}{x^{\beta+1}}, \quad x > \alpha,$$

where  $\alpha > 0$  and  $\beta > 0$ . If  $X$  has a Pareto distribution, then

$$E[X] = \frac{\beta\alpha}{\beta - 1}, \quad \beta > 1$$

$$\text{Var}[X] = \frac{\beta\alpha^2}{(\beta - 1)^2(\beta - 2)}, \quad \beta > 2,$$

where for  $0 < \beta \leq 1$  the expected value is infinite and for  $0 < \beta \leq 2$  the variance is infinite.

**14. Extreme value (or Gumbel) Distribution.** The PDF for the extreme value distribution is

$$f(x) = e^{-e^{-\frac{x-\alpha}{\beta}}} \frac{1}{\beta} e^{-\frac{x-\alpha}{\beta}},$$

where  $-\infty < \alpha < \infty$  and  $\beta > 0$ . If  $X$  has an extreme value distribution, then

$$E[X] = \alpha + \beta\gamma$$

$$\text{Var}[X] = \frac{\pi^2\beta^2}{6}.$$

Here  $\gamma \cong .577216$  is *Euler's constant*.<sup>14</sup>

**15. Chi-square Distribution.** The PDF for the chi-square distribution is

$$f(x) = \frac{1}{\Gamma(k/2)} \left(\frac{1}{2}\right)^{k/2} x^{k/2-1} e^{-\frac{1}{2}x}, \quad x > 0,$$

where  $k = 1, 2, \dots$  is called the *number of degrees of freedom*. If  $X$  has a chi-square distribution, then

$$E[X] = k$$

$$\text{Var}[X] = 2k.$$

The chi-square distribution is often the distribution satisfied by a test statistic in hypothesis testing. It is precisely the distribution of the sum of the squares of  $k$  independent standard normal random variables,

$$X = N_1^2 + N_2^2 + \dots + N_k^2.$$

**16. t distribution.** The PDF for the t distribution is

$$f(x) = \frac{\Gamma(\frac{k+1}{2})}{\Gamma(\frac{k}{2})} \frac{1}{\sqrt{k\pi}} \frac{1}{(1 + x^2/k)^{(k+1)/2}},$$

---

<sup>14</sup>Not to be confused with *Euler's number*,  $e \cong 2.7183$ .

where  $k > 0$ . If  $X$  has a t distribution, then

$$E[X] = 0, \quad k > 1$$

$$\text{Var}[X] = \frac{k}{k-2}, \quad k > 2.$$

The t distribution is the distribution for

$$X = \frac{N}{\sqrt{\chi/k}},$$

where  $N$  is a standard normal random variable, and  $\chi$  is a chi-square random variable with  $k$  degrees of freedom.

**17. F distribution.** The PDF for the F distribution is

$$f(x) = \frac{\Gamma(\frac{r_1+r_2}{2})}{\Gamma(\frac{r_1}{2})\Gamma(\frac{r_2}{2})} r_1^{r_1/2} r_2^{r_2/2} \frac{x^{r_1/r_2-1}}{(r_2 + r_1 x)^{(r_1+r_2)/2}}, \quad x > 0,$$

where  $m = 1, 2, \dots$  and  $n = 1, 2, \dots$ . If  $X$  has an  $F$  distribution, then

$$E[X] = \frac{r_2}{r_2 - 2}, \quad r_2 > 2$$

$$\text{Var}[X] = \frac{2r_2^2(r_1 + r_2 - 2)}{r_1(r_2 - 2)^2(r_2 - 4)}, \quad r_2 > 4.$$

The F distribution is the distribution for

$$X = \frac{\chi_1}{r_1} \frac{\chi_2}{r_2},$$

where  $\chi_1$  and  $\chi_2$  are independent chi-square random variables with  $r_1$  and  $r_2$  degrees of freedom respectively.

## 7.7 Joint Probability Density Functions

**Definition.** Given two random variables  $X$  and  $Y$ , we define the joint cumulative probability distribution function as

$$F_{X,Y}(x, y) = P(X \leq x, Y \leq y).$$

**Definition 7.4.** We will say that that a function  $f(x, y)$  is the joint probability density function for a pair of continuous random variables  $X$  and  $Y$  provided  $f$  has the following properties:

- $f(x, y) \geq 0$  for all  $x, y \in \mathbb{R}$ ;
- $\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f(x, y) dx dy = 1$ ;

- For any rectangle  $R = [a, b] \times [c, d] \subset \mathbb{R} \times \mathbb{R}$ ,  $a < b$ ,  $c < d$ ,

$$P((X, Y) \in R) = \int \int_R f(x, y) dx dy,$$

where it's assumed in this last item that the integral makes sense for all such rectangles.

**Remark 7.2.** Similarly as in the case of PDFs for single random variables, this is not the most general definition of either a pair of jointly distributed continuous random variables or a joint probability density function. For the integral condition, it's again appropriate to think of  $f$  as being Riemann integrable on each such rectangle  $R$ , or more generally just Lebesgue measurable on all such rectangles. We also note that whenever  $F$  is twice differentiable, we have the relation

$$f(x, y) = \frac{\partial^2}{\partial x \partial y} F(x, y).$$

## 8 Maximum Likelihood Estimators

Once we decide on a probability density function that we expect will reasonably fit our data, we must use our data to determine values for the parameters of the distribution. One method for finding such parameter values is that of maximum likelihood estimation.

### 8.1 Maximum Likelihood Estimation for Discrete Random Variables

**Example 8.1.** Consider a coin, known to be unfair, with probability of heads either  $p = .6$  or  $p = .4$ , and suppose we would like to use experimental data to determine which is correct. That is, we are trying to estimate the value of the parameter  $p$ .

First, consider the case in which we are only given one flip on which to base our decision. Let  $X$  be a random variable denoting the number of heads that turn up in a given flip, and let  $f(x; p) = P(X = x|p)$  be the associated probability mass function. We have two possibilities for  $f$ ,

$$f(x; .6) = \begin{cases} .4 & X = 0 \\ .6 & X = 1, \end{cases} \quad f(x; .4) = \begin{cases} .6 & X = 0 \\ .4 & X = 1. \end{cases}$$

It will be useful to think in terms of a table of possible outcomes Table 8.1.

PDF/Number of Heads that turn up in experiment	$x = 0$	$x = 1$
$f(x; .6)$	.4	.6
$f(x; .4)$	.6	.4

Table 8.1: Analysis of an unfair coin with a single flip.

Clearly, the only conclusion we can make from a single flip is that if  $X = 1$  (the coin turns up heads), we take  $p = .6$ , while if  $X = 0$ , we take  $p = .4$ .

Next, suppose we have two flips  $X_1$  and  $X_2$ , so that the probability mass function becomes  $f(x_1; p)f(x_2; p)$ , where  $f$  is the PDF for a single flip. Proceeding as above, we have Table 8.2. If the outcome of  $X_1$  and  $X_2$  are both 0, then we have two options. Either  $p = .6$  and this outcome had probability  $.4^2 = .16$ , or  $p = .4$  and this outcome had probability  $.6^2 = .36$ . Since  $x_1 = 0, x_2 = 0$  is the outcome actually observed in this case, we might expect that it is the most likely of the two possibilities to observe, and we choose  $p = .4$ . We can argue similarly for the other possibilities, taking  $p = .6$  if the outcomes of  $X_1$  and  $X_2$  are both 1, and asserting no difference if one of the outcomes is 0 while the other is 1. In all of these cases, the upshot is that we choose the parameter value  $p$  that maximizes the value of the relevant probability mass function ( $f(x; p)$  for a single observation,  $f(x_1; p)f(x_2; p)$  for two observations). This process of maximizing the probability mass function is what we refer to as maximum likelihood estimation.

PDF/Number of heads	$x_1 = 0, x_2 = 0$	$x_1 = 0, x_2 = 1$	$x_1 = 1, x_2 = 0$	$x_1 = 1, x_2 = 1$
$f(x_1; .6)f(x_2; .6)$	$.4^2$	$.4 \cdot .6$	$.6 \cdot .4$	$.6^2$
$f(x_1; .4)f(x_2; .4)$	$.6^2$	$.6 \cdot .4$	$.4 \cdot .6$	$.4^2$

Table 8.2: Analysis of an unfair coin two flips.

## 8.2 Maximum Likelihood Estimation for Continuous Random Variables

The main observation we take from our discussion of maximum likelihood estimators in the case of discrete random variables is that we should choose  $p$  so that the relevant probability mass function is maximized. In this way, we observe that finding a maximum likelihood estimator is a maximization problem, and in the continuous case, we will be able to use methods from calculus.

**Example 8.2.** Suppose an experimental set of measurements  $x_1, x_2, \dots, x_n$  appears to have arisen from an exponential distribution. Determine an MLE for the parameter  $a$ .

As in Example 8.1, we first consider the case of a single measurement,  $x_1$ . The PDF for the exponential distribution is

$$f(x; a) = ae^{-ax}, \quad x > 0,$$

and so we search for the value of  $a$  that maximizes

$$L(a) = f(x_1; a),$$

which is called the *likelihood function*. For this, we keep in mind that our rationale here is precisely as it was in the discrete case: since  $x_1$  is the observation we get, we want to choose  $a$  so as to make this as likely as possible—the assumption being that in experiments we most often see the most likely events. Here, we have

$$L'(a) = e^{-ax_1} - ax_1e^{-ax_1} = 0 \implies e^{-ax_1}(1 - ax_1) = 0 \implies a = \frac{1}{x_1}.$$

In the case of  $n$  measurements  $x_1, x_2, \dots, x_n$ , the likelihood function becomes the joint PDF

$$\begin{aligned} L(a) &= f(x_1; a)f(x_2; a) \cdots f(x_n; a) \\ &= \prod_{k=1}^n f(x_k; a) \\ &= \prod_{k=1}^n a e^{-ax_k} \\ &= a^n e^{-a \sum_{k=1}^n x_k}. \end{aligned}$$

In this case,

$$\frac{\partial L}{\partial a} = na^{n-1}e^{-a \sum_{k=1}^n x_k} - a^n \left( \sum_{k=1}^n x_k \right) e^{-a \sum_{k=1}^n x_k},$$

so that we have,

$$a = \frac{n}{\sum_{k=1}^n x_k}.$$

(We observe that since we are maximizing  $L(a)$  on the domain  $0 \leq a < \infty$ , we need only check that  $L(0) = 0$  and  $\lim_{a \rightarrow \infty} L(a) = 0$  to see that this is indeed a maximum.) In order to simplify calculations of this last type, we often define the *log-likelihood* function,

$$L^*(a) = \ln(L(a)),$$

simply the natural logarithm of the likelihood function. The advantage in this is that since natural logarithm is monotonic,  $L$  and  $L^*$  are maximized by the same value of  $a$ , and due to the rule of logarithms, if  $L(a)$  is a product,  $L^*(a)$  will be a sum. Also, in the case of a large number of data point,  $L$  can become quite large. In the current example,

$$L^*(a) = \ln(a^n e^{-a \sum_{k=1}^n x_k}) = \ln(a^n) + \ln(e^{-a \sum_{k=1}^n x_k}) = n \ln a - a \sum_{k=1}^n x_k,$$

so that

$$\frac{\partial L^*}{\partial a} = \frac{n}{a} - \sum_{k=1}^n x_k = 0,$$

which gives the same result as obtained above. △

More generally, we can compute maximum likelihood estimators with the aid of MATLAB.

**Example 8.3.** Given a set of 100 data points that appear to arise from a process that follows a Weibull distribution, determine maximum likelihood estimators for  $\lambda$  and  $\beta$ .

The data for this example can be computed with the M-file *weibsim.m*: (For a discussion of simulating random variables, see Section 9.)

```

%WEIBSIM: MATLAB script M-file that simulates 100
%Weibull random variables
lam = .5; bet = 2;
for k=1:100
X(k) = (1/lam)*(-log(1-rand))^(1/bet);
end

```

We first recall that the PDF for the Weibull distribution is

$$f(x) = \begin{cases} \lambda^\beta \beta x^{\beta-1} e^{-(\lambda x)^\beta} & x > 0 \\ 0 & x < 0, \end{cases}$$

and so the likelihood function is

$$L(\lambda, \beta) = \prod_{k=1}^n \lambda^\beta \beta x_k^{\beta-1} e^{-(\lambda x_k)^\beta}.$$

In theory, of course, we can find values of  $\lambda$  and  $\beta$  by solving the system of equations

$$\begin{aligned} \frac{\partial L}{\partial \lambda} &= 0 \\ \frac{\partial L}{\partial \beta} &= 0. \end{aligned}$$

In practice, however, we employ MATLAB. First, we record the likelihood function in a function M-file, which takes values  $\lambda$ ,  $\beta$ , and  $\{x_k\}_{k=1}^n$ , and returns values of the log-likelihood function  $\ln(L)$  (see *weiblike.m*). (In this case, we use the log-likelihood function because of how large this product of 100 points becomes.)

```

function value = weiblike(p, X);
%WEIBLIKE: MATLAB function M-file that compute the likelihood function
%for a Weibull distrubution with data vector D
%Note: D is passed as a parameter
%p(1) = lambda, p(2) = beta
value = 0;
for k=1:length(X)
value=value+log(p(1)^(p(2))*p(2)*X(k)^(p(2)-1)*exp(-(p(1)*X(k))^(p(2))));
end
value = -value;

```

We observe that since MATLAB's optimization routines are for minimization, we multiply the function by a negative sign for maximization. In order to find the optimal values and plot our fit along with a histogram, we use *weibfit.m*.

```

function weibfit(X)
%WEIBFIT: MATLAB function M-file that fits data to a Weibull
%distribution, and checks fit by plotting a scaled PDF along
%with a histogram.

```



```

hold off;
guess = [1 1];
options=optimset('MaxFunEvals',10000);
[p, LL]=fminsearch(@weiblike,guess,options,X)
%Plotting
hist(X,12)
[n,c]=hist(X,12);
hold on;
x = linspace(min(X),max(X),50);
fweib = sum(n)*(c(2)-c(1))*p(1)^(p(2))*p(2)*x.^(p(2)-1).*exp(-(p(1)*x).^(p(2)));
plot(x,fweib,'r')

```

We obtain the fit given in Figure 8.1, and the parameter values are  $\lambda = .5029$  and  $\beta = 2.2368$ .

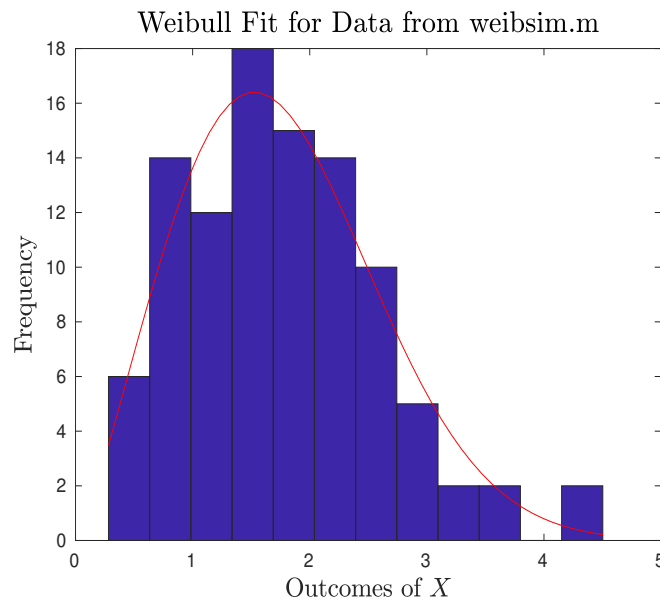


Figure 8.1: MLE PDF fit for Weibull distribution.

## 9 Simulating a Random Process

Sometimes the best way to determine how a certain phenomenon will play out is to simulate it several times and simply watch what happens. Generally, this is referred to as the *Monte Carlo* method, after a famous casino by that name in Monaco (a French province on the Mediterranean).<sup>15</sup>

<sup>15</sup>Apparently, this name was given to the method by Nicholas Metropolis, a researcher at Los Alamos National Laboratory (located in New Mexico) in the 1940s, and was used to describe simulations they were running of nuclear explosions. The first appearance of the method appears to be in the 1949 paper, “The Monte Carlo Method,” published by Metropolis and Stanislaw Ulam in the Journal of the American Statistical Association.

**Example 9.1.** What is the expected number of flips of a fair coin until it turns up heads?

At first glance, this might look like a difficult problem to study analytically. The problem is that if we begin computing the expected value from the definition, we get an infinite series,

$$E[N] = 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} + 3 \cdot \frac{1}{8} + \dots$$

That is, the probability that it takes one flip is  $\frac{1}{2}$ ; the probability that it takes two flips is  $\frac{1}{4}$  etc. In fact, we have already developed a method for analyzing problems like this in Section 3, but for now let's suppose we want to take an alternative approach. One method is to simply pull out a coin and begin flipping it, counting how many flips it takes for it to land heads. This should lead to some sort of list, say 2, 4, 1, 3, 2 for five trials. The average of these should give us an approximation for the expected value:

$$E[N] \cong \frac{2 + 4 + 1 + 3 + 2}{5} = \frac{12}{5}.$$

The more trials we run, the better our approximation should be.

In general, we would like to carry out such simulations computationally. We accomplish this with *pseudo random numbers*, which behave randomly as long as we don't watch them too carefully.<sup>16</sup> Our fundamental random variable generator from MATLAB will be the built-in function *rand*, which creates a real number, fifteen digits long, uniformly distributed on the interval [0,1]. (The fact that this number has a finite length means that it is not really a continuous random variable, but with fifteen digits it will approximate the outcome of a continuous random variable sufficiently well for all calculations carried out in these notes.) In the following MATLAB code, we take *rand*<=.5 to correspond with the coin landing on heads and *rand*>.5 to correspond with it landing tails.

```
function ev = flips(n)
%FLIPS: MATLAB function M-file that simulates
%flipping a coin until it turns up heads. The
%input, n, is number of trials, and the
%output, ev, is expected value.
for k=1:n
m=1; %m counts the number of flips
while rand > .5 %while tails
m=m+1;
end
R(k)=m;
end
ev=mean(R);
```

We now compute as follows in the MATLAB Command Window.

---

<sup>16</sup>There is an enormous amount of literature regarding pseudo random numbers, which we will ignore entirely. For our purposes, the random variables MATLAB creates will be sufficient. If, however, you find yourself doing serious simulation (i.e., getting paid for it) you should *at least* understand the generator you are using.

```

>>flips(10)
ans =
1.6000
>>flips(100)
ans =
1.8700
>>flips(1000)
ans =
2.0430
>>flips(10000)
ans =
1.9958

```

Observe that as we take more trials, the mean seems to be converging to 2. △

## 9.1 Simulating Uniform Random Variables

As mentioned above, the MATLAB built-in function *rand* creates pseudo random numbers uniformly distributed on the interval  $[0, 1]$ . In order to develop a new random variable,  $U$ , uniformly distributed on the interval  $[a, b]$ , we need only use  $U=a+(b-a)*rand$ .

## 9.2 Simulating Discrete Random Variables

We can typically build discrete random variables out of uniform random variables by either conditioning (through *if* or *while* statements) or rounding. Suppose we want to simulate a random variable that is 0 with probability  $1/2$  and 1 with probability  $1/2$ . We can either condition,

```

if rand < .5
X = 0;
else
X = 1;
end

```

or round

```

X=round(rand)

```

Notice that in either case we ignore the subtle problem that the probability that  $rand < .5$  is slightly smaller than the complementary probability that  $rand \geq .5$ , which includes the possibility of equality. Keep in mind here that MATLAB computes *rand* to fifteen decimal places of accuracy, and so the probability that *rand* is precisely  $.5$  is roughly  $10^{-14}$ .

As another example, suppose we want to simulate the role of a fair die. In this case, our *if* statement would grow to some length, but we can equivalently use the single line,

```

R=round(6*rand+.5)

```

(Here, the addition of .5 simply insures that we never get a roll of 0.)

Another option, similar to using *round*, is the use of MATLAB's function *ceil* (think *ceiling*), which rounds numbers to the nearest larger integer. See also *floor*.

Finally, MATLAB's built-in M-file *randi.m* simulates any number of random integers.

### 9.3 Simulating Gaussian Random Variables

MATLAB also has a built-in Gaussian random number generator, *randn*, which creates pseudo random numbers from a Gaussian distribution with mean 0 and variance 1 (such a distribution is also referred to as the *standard normal* distribution). In order to see how we can generate more general Gaussian random numbers, we let  $N$  denote a standard normal random variable; that is,  $E[N] = 0$  and  $\text{Var}[N] = E[(N - E[N])^2] = E[N^2] = 1$ . Introducing the new random variable  $X = \mu + \sigma N$ , we have

$$E[X] = E[\mu + \sigma N] = E[\mu] + \sigma E[N] = \mu,$$

and

$$\begin{aligned}\text{Var}[X] &= E[(X - E[X])^2] = E[X^2 - 2\mu X + \mu^2] = E[(\mu + \sigma N)^2 - 2\mu(\mu + \sigma N) + \mu^2] \\ &= E[\mu^2 + 2\sigma\mu N + \sigma^2 N^2 - 2\mu^2 - 2\mu\sigma N + \mu^2] = \sigma^2 E[N^2] = \sigma^2,\end{aligned}$$

from which we *suspect* that  $X$  is a Gaussian random variable with mean  $\mu$  and standard deviation  $\sigma$ . In order to create pseudo random numbers from such a distribution in MATLAB, with mean *mu* and standard deviation *sigma*, we simply use  $X = \text{mu} + \text{sigma} * \text{randn}$ .

In the previous discussion, we have not actually proven that  $X$  is a Gaussian random variable, only that it has the correct expected value and variance. In order to prove that it is Gaussian distributed, we must show that it has the Gaussian PDF. In order to do this, we will compute its CDF and take a derivative. We compute

$$\begin{aligned}F(x) &= P(X \leq x) = P(\mu + \sigma N \leq x) = P(N \leq \frac{x - \mu}{\sigma}) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{x - \mu}{\sigma}} e^{-\frac{y^2}{2}} dy.\end{aligned}$$

We have, then, according to the fundamental theorem of calculus

$$f(x) = F'(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x - \mu)^2}{2\sigma^2}},$$

which is indeed the PDF for a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ .

### 9.4 Simulating More General Random Variables

In order to simulate general random variables, we will require two theorems.

**Theorem 9.1.** *Suppose the random variable  $X$  has a cumulative distribution function  $F(x)$ , where  $F(x)$  is continuous and strictly increasing whenever it is not 0 or 1. Then  $X \stackrel{d}{=}$*

$F^{-1}(Y)$ , where  $Y$  is uniformly distributed on  $[0, 1]$  and by  $\stackrel{d}{=}$  we mean equal in distribution: that each random variable has the same distribution.

**Proof.** First, we observe that for  $y \in [0, 1]$  and  $Y$  uniformly distributed on  $[0, 1]$ , we have  $P(Y \leq y) = y$ . Next, we note that our assumptions of continuity and monotonicity on  $F(x)$  require it to behave somewhat like the example cumulative distribution function sketched in Figure 9.1.

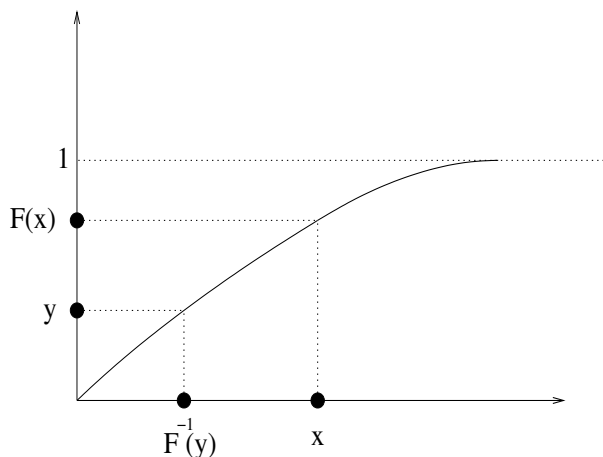


Figure 9.1:  $F(x)$  continuous and strictly increasing.

We have, then, that the cumulative distribution function for  $X = F^{-1}(Y)$  is given by

$$F_{F^{-1}(Y)}(x) = P(F^{-1}(Y) \leq x) = P(Y \leq F(x)) = F(x),$$

where the first equality follows from the definition of cumulative distribution function, the second follows from continuity and monotonicity, and the third follows from our first observation of the proof.  $\square$

**Example 9.2.** Assuming  $Y$  is a uniformly distributed random variable on  $[0, 1]$ , develop a random variable  $X$  in terms of  $Y$  that satisfies the exponential distribution.

First, we compute  $F(x)$  for the exponential distribution by integrating over the probability density function,

$$F(x) = \int_0^x ae^{-ay} dy = -e^{-ay} \Big|_0^x = 1 - e^{-ax}, \quad x \geq 0.$$

Clearly,  $F(x)$  satisfies the conditions of Theorem 7.1, so all that remains is to find  $F^{-1}$ . We write  $y = 1 - e^{-ax}$  and solve for  $x$  to find  $x = -\frac{1}{a} \log(1 - y)$ , or in terms of  $X$  and  $Y$ ,  $X = -\frac{1}{a} \log(1 - Y)$ .  $\triangle$

#### 9.4.1 The Rejection Method.

A more general method for simulating random variables is the *rejection* method. Suppose we can simulate random variables associated with some probability density function  $g(x)$ , and would like to simulate random variables from a second probability density function  $f(x)$ . The rejection method follows the steps outlined below.

1. Let  $c$  be a constant so that

$$\frac{f(x)}{g(x)} \leq c \quad \text{for all } x \in \mathbb{R},$$

where for efficiency  $c$  is to be chosen as small as possible.

2. Simulate both a random variable  $Y$  with density  $g(y)$  and a random variable  $U$  uniformly distributed on  $[0, 1]$ .

3. If  $U \leq \frac{f(Y)}{cg(Y)}$ , set  $X = Y$ . Otherwise, repeat Step 2.

**Theorem 9.2.** *The random variable  $X$  created by the rejection method has probability density function  $f(x)$ .*

**Proof.** Let  $X$  be the random variable created by the rejection method, and compute its associated cumulative distribution function,

$$F_X(x) = P(X \leq x) = P(Y \leq x | U \leq \frac{f(Y)}{cg(Y)}) = \frac{P(\{Y \leq x\} \cap \{U \leq \frac{f(Y)}{cg(Y)}\})}{P(U \leq \frac{f(Y)}{cg(Y)})}.$$

Since  $Y$  and  $U$  are independent random variables, the joint probability density function of  $Y$  and  $U$  is

$$p(y, u) = g(y)f_U(u),$$

where  $f_U(u)$  is the probability density function for  $U$ . We observe that

$$\begin{aligned} P(\{Y \leq x\} \cap \{U \leq \frac{f(Y)}{cg(Y)}\}) &= \int_{-\infty}^x \int_0^{\frac{f(y)}{cg(y)}} f_U(u)g(y)dudy \\ &= \int_{-\infty}^x g(y) \int_0^{\frac{f(y)}{cg(y)}} dudy \\ &= \int_{-\infty}^x \frac{f(y)}{c} dy. \end{aligned}$$

We have, then,

$$P(X \leq x) = \frac{1}{cP(U \leq \frac{f(Y)}{cg(Y)})} \int_{-\infty}^x f(y)dy.$$

Taking the limit as  $x \rightarrow \infty$ , we see that  $cP(U \leq \frac{f(Y)}{cg(Y)}) = 1$ , and consequently  $F_X(x) = \int_{-\infty}^x f(y)dy$ .  $\square$

**Example 9.3.** Develop a MATLAB program that simulates a beta random variable (in the case that the beta PDF is bounded).

For simplicity, we take our known probability density function  $g(y)$  to be uniformly distributed on the interval  $[0, 1]$ ; i.e.,

$$g(y) = \begin{cases} 1, & 0 \leq y \leq 1, \\ 0, & \text{otherwise.} \end{cases}$$

For  $c$  depending on the values of  $a$  and  $b$ , we simulate both  $Y$  and  $U$  as (independent) uniformly distributed random variables on  $[0, 1]$  and check of  $U \leq \frac{f(Y)}{c}$ , where  $f(Y)$  denotes the beta probability density function evaluated at  $Y$ . If  $U \leq \frac{f(Y)}{c}$ , we set  $X = Y$ , otherwise we repeat Step 2. This process is carried out in the MATLAB function M-file *ranbeta.m*.

```
function b = ranbeta(a,b,c);
%RANBETA: MATLAB function M-file for simulating a
%random variable with a beta distribution. The value of c
%must be greater than the maximum of the associated beta
%PDF, though not by much, or too many iterations will
%be required. This file employs the rejection method with
%comparison PDF g uniform on [0,1]; i.e., identically 1.
m = 0;
while m<1
var = rand; %Simulates Y
f = (1/beta(a,b))*var.^(a - 1).*(1 - var).^(b - 1);
if rand <= f/c %rand is simulating U
b = var;
m = 1;
end
end
```

As an example implementation, we will take  $a = 2$  and  $b = 4$ . In this case,

$$f(x) = 20x(1 - x)^3.$$

In order to select an appropriate value for  $c$ , we can find the maximum value of  $f$ . Setting  $f'(x) = 0$ , we find that the maximum occurs at  $x = 1/4$ , which gives  $|f(x)| \leq 2.1094$ . We can choose  $c = 2.2$ . We will simulate data and plot a histogram of the data along with a scaled pdf with the MATLAB M-file *betafit.m*.

```
%BETAFIT: MATLAB script M-file that uses betasim.m to simulate
%a beta random variable and tests the result by plotting
%a scaled PDE along with a histogram
a = 2; b = 4; c = 2.2;
for k=1:1000
B(k) = ranbeta(a,b,c);
end
hist(B,12)
[n,c]=hist(B,12);
hold on;
x = linspace(0,1,100);
fbeta = sum(n)*(c(2)-c(1))/(beta(a,b))*x.^(a-1).*(1-x).^(b-1);
plot(x,fbeta,'r')
```

The resulting figure is Figure 9.2.

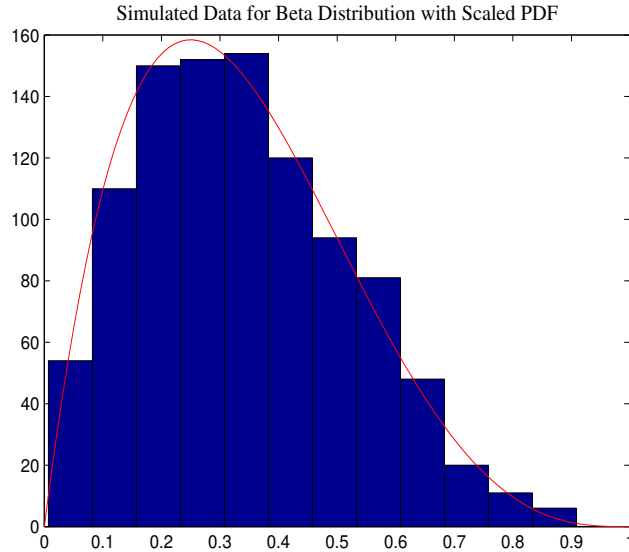


Figure 9.2: Plot of simulated beta distributed data along with scaled pdf.

## 9.5 Application to Queuing Theory

In this section, we consider an application to the study of queueing theory<sup>17</sup>. Suppose we want to simulate various situations in which customers randomly arrive at some service station. The following script M-file, *queue1.m*, simulates the situation in which the customer arrival times are exponentially distributed and the service times are fixed at exactly .5 minutes per customer (or  $\mu = 2$  customers/minute).

```
%QUEUE1: Script file for simulating customers
%arriving at a single queue.
S=.5; %Time to service customers (1/mu)
m=.5; %Mean time between customer arrivals (1/lambda)
Q=0; %Time to service all customers remaining in queue (I.e., time until
% an arriving customer is served.)
queuwait=0; %Total time customers spent waiting in queue
systemwait=0; %Total time customers spent waiting in system
N0=0; %Number of times arriving customer finds no one in queue
N=1000; %Number of customers to watch
% The simulation
for k=1:N %Watch N customers
T=-m*log(1-rand); %Arrival of customer C (exponential)
Q=max(Q-T,0); %Time T has elapsed since time Q was set
if Q==0
```

<sup>17</sup>Sometimes spelled *Queueing* Theory, as for example in Sheldon Ross's influential textbook *Introduction to Probability Models*, Academic Press 1989. In adopting the convention *queueing*, I'm following Bryan A. Garner's *A Dictionary of Modern American Usage*, Oxford University Press, 1998.



```

N0=N0+1;
end
queuwait=queuwait+Q; %Total time spent waiting in queue
Q = Q + S; %Time until customer C leaves system
systemwait=systemwait+Q; %Total time spent waiting in system
end
Wq=queuwait/N
W=systemwait/N
P0=N0/N

```

## 9.6 Limit Theorems

When proceeding by simulation we typically make the following pair of assumptions:

1. If we take a large number of observations  $X_1, X_2, \dots, X_n$  ( $n$  large) of the same process  $X$ , the average of these observations will be a good approximation for the average of the process:

$$E[X] \approx \frac{1}{n} \sum_{k=1}^n X_k.$$

2. If we define a random variable  $Y$  as the sum of  $n$  repeatable observations,

$$Y = \sum_{k=1}^n X_k,$$

then  $Y$  will be approximately Gaussian.

In this section we will investigate conditions under which these statements are justified. Results of this kind are typically referred to as Limit Theorems: those of Type 1 are “laws of large numbers,” while those of type 2 are “central limit theorems.” We will start with a lemma that’s named for the Russian mathematician Andrei Markov (1856-1922), though it was apparently first published by his doctoral advisor Pafnuty Chebyshev (1821-1894).

**Lemma 9.1.** (Markov’s Inequality) *If  $X$  is a random variable that takes only non-negative values, then for any  $a > 0$*

$$P(X \geq a) \leq \frac{E[X]}{a}.$$

**Proof.** For  $a > 0$  set

$$I_{[a,\infty)}(x) := \begin{cases} 1 & x \geq a \\ 0 & x < a \end{cases},$$

which is typically referred to as the *indicator function* for the set  $[a, \infty)$ . Let  $Y := I_{[a,\infty)}(X)$  be a new random variable, and observe that  $Y \leq \frac{X}{a}$ . I.e., if  $0 \leq X < a$ , then  $Y = 0 \leq \frac{X}{a}$ , while if  $X \geq a$ , then  $Y = 1 \leq \frac{X}{a}$ . Then

$$E[Y] = 0 \cdot P(X < a) + 1 \cdot P(X \geq a),$$

so that

$$P(X \geq a) = E[Y] \leq E\left[\frac{X}{a}\right] = \frac{E[X]}{a}.$$

which is the claim.  $\square$

The advantage of an inequality like this is that it allows us to obtain information about certain probabilities without knowing how the random variables involved are distributed.

**Example 9.4.** Suppose that a certain student's exam average in a class is 75 out of 100, and that each exam is equally difficult. Find an upper bound on the probability that this student will make a 90 on the next exam.

Let  $X$  be a random variable corresponding with the score a student will make on any exam. We have  $E[X] = 75$ , so using Markov's inequality, we find

$$P(X \geq 90) \leq \frac{E[X]}{90} = \frac{75}{90} = \frac{5}{6}.$$

(Note that this calculation hinges on the fairly implausible assumption that the student's expected exam score is precisely  $E[X] = 75$ . On the other hand, we do not have to know anything about how  $X$  is distributed.)  $\triangle$

**Lemma 9.2.** (Chebyshev's Inequality) *If  $X$  is a random variable with finite mean  $\mu$  and standard deviation  $\sigma$  (not necessarily finite), then for any  $k > 0$*

$$P(|X - \mu| \geq k) \leq \frac{\sigma^2}{k^2}.$$

**Proof.** We can prove this by applying Markov's inequality to the non-negative random variable  $Z = (X - \mu)^2$ . That is,

$$P((X - \mu)^2 \geq k^2) \leq \frac{E[(X - \mu)^2]}{k^2},$$

which is equivalent to the claim.  $\square$

**Example 9.5.** For the student discussed in Example 9.4 suppose the standard deviation for  $X$  is 5. What is the probability the student's grade on the final will be between 60 and 90?

Using Chebyshev's inequality, we have

$$P(|X - 75| \geq 15) \leq \frac{25}{15^2} = \frac{1}{9},$$

and so

$$P(|X - 75| < 15) \geq \frac{8}{9}.$$

It should be fairly clear that both Markov's inequality and Chebyshev's inequality can give crude results. For comparison, in the next example, we consider the case in which  $X$  is known to be Gaussian.

**Example 9.6.** Suppose the random variable  $X$  from Examples 9.4 and 9.5 is known to be Gaussian. Compute  $P(X \geq 90)$ .

In this case  $\mu = 75$  and  $\sigma^2 = 25$ , so the Gaussian probability density function is

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

and consequently

$$P(X \geq 90) = \int_{90}^{\infty} \frac{1}{\sqrt{50\pi}} e^{-\frac{(x-75)^2}{50}} dx = .0013,$$

or .13%. △

**Theorem 9.3.** (The Weak Law of Large Numbers) *Let  $\{X_j\}_{j=1}^{\infty}$  denote a sequence of independent and identically distributed (i.i.d.) random variables with  $E[|X_1|]$  finite (so  $E[|X_j|]$  finite for all  $j \in \mathbb{N}$ ). Then  $\mu = E[X_1]$  is finite, and given any  $\epsilon > 0$*

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{1}{n} \sum_{j=1}^n X_j - \mu\right| \geq \epsilon\right) = 0.$$

**Remark 9.1.** *The weak law of large numbers was first proven by the Swiss mathematician Jacob Bernoulli (1654–1705) for the special case of Bernoulli random variables, and in the form stated here by the Russian Mathematician Aleksandr Yakovlevich Khintchine (1894–1959). This type of convergence is called convergence in probability, and is sometimes referred to as weak convergence.*

**Proof.** Though the theorem is true regardless of whether or not the variance associated with these random variables is finite, we prove it only for the case of finite variance. Under this assumption, we'll write

$$E[X_j] = \mu \quad \text{and} \quad \text{Var}[X_j] = \sigma^2,$$

for all  $j \in \mathbb{N}$ . Our approach will be to apply Chebyshev's Inequality to the random variable  $Y := \frac{1}{n} \sum_{j=1}^n X_j$ .

First, we notice that

$$E[Y] = \frac{1}{n} E\left[\sum_{j=1}^n X_j\right] = \frac{1}{n} \sum_{j=1}^n E[X_j] = \frac{1}{n} (\mu + \mu + \cdots + \mu) = \mu,$$

and also

$$\text{Var}[Y] = \text{Var}\left[\frac{1}{n} \sum_{j=1}^n X_j\right] = \frac{1}{n^2} \sum_{j=1}^n \text{Var}[X_j] = \frac{\sigma^2}{n}.$$

Chebyshev's Inequality then asserts

$$P(|Y - \mu| \geq \epsilon) \leq \frac{\frac{\sigma^2}{n}}{\epsilon^2} = \frac{\sigma^2/n}{\epsilon^2}. \tag{9.1}$$

Now we simply take the limit as  $n \rightarrow \infty$ . □

In practical applications, inequality (9.1) is often more useful than the full theorem.

**Example 9.7.** Suppose a trick coin is to be flipped  $n$  times, and we would like to determine from this experiment the probability that it will land heads. Determine the number of experiments  $n$  necessary to provide 95% certainty that the error is less than .01.

We can proceed by defining a random variable  $X$  as follows:

$$X = \begin{cases} 1 & \text{if the coin lands heads} \\ 0 & \text{if the coin lands tails,} \end{cases}$$

for which

$$E[X] = 1 \cdot P(H) + 0 \cdot P(T),$$

so that if we set  $p = P(H)$  then  $p = \mu$ . Letting now  $X_j$  denote the random variable corresponding with the outcome of the  $j^{\text{th}}$  flip, we set

$$\hat{p} := \frac{1}{n} \sum_{k=1}^n X_k.$$

Here,  $\hat{p}$  is a random variable called an *estimator* for the probability  $p$ . Notice that similarly as in the proof of Theorem 9.3,

$$E[\hat{p}] = \mu = p \quad \text{and} \quad \text{Var}[\hat{p}] = \frac{\sigma^2}{n}.$$

Now, according to the weak law of large numbers (in particular, inequality (9.1)),

$$P(|\hat{p} - p| \geq \epsilon) \leq \frac{\sigma^2}{n\epsilon^2}.$$

If we want two decimal places of accuracy in our approximation of  $p$ , we take  $\epsilon = .01$ , and for 95% certainty we want

$$\frac{\sigma^2}{n(.01)^2} \leq .05.$$

Here,

$$\sigma^2 = \text{Var}[X] = E[X^2] - E[X]^2,$$

but in this case  $X = X^2$ , so in fact  $\sigma^2 = p - p^2$ . We must have  $p \in [0, 1]$ , so a simple optimization shows that  $\sigma^2 \leq \frac{1}{4}$ . We can write

$$\frac{\sigma^2}{n(.01)^2} \leq \frac{1}{4n(.01)^2} \leq .05,$$

where we would like to force the second inequality to be true by choosing  $n$  sufficiently large. In particular, we can take

$$4n(.01)^2 \geq \frac{1}{.05} \implies n \geq \frac{200000}{4} = 50000.$$

We can conclude that if we conduct 50,000 experiments, then we will have  $|\bar{p} - p| < .01$  with probability at least 95%.  $\triangle$

**Theorem 9.4.** (Strong Law of Large Numbers) *Let  $\{X_j\}_{j=1}^\infty$  be a sequence of i.i.d. random variables with  $E[|X_1|]$  finite, and set  $\mu = E[X_1]$ . Then*

$$P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n X_j = \mu\right) = 1.$$

**Remark 9.2.** *We recall that the Weak Law asserts*

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{1}{n} \sum_{j=1}^n X_j - \mu\right| \geq \epsilon\right) = 0.$$

*This says that if a large enough number of experiments is carried out, then the probability that  $\frac{1}{n} \sum_{j=1}^n X_j$  differs from  $\mu$  by more than  $\epsilon$  will be small. But a difference of size  $\epsilon$  is allowed even in the limit. The Strong Law asserts, rather, that if we could, in principle, conduct an infinite number of experiments, we would certainly obtain the mean (i.e., with probability 1). The type of convergence in the Strong Law is called almost sure convergence or strong convergence. It's possible to prove in general that strong convergence implies weak convergence, but not the converse, justifying the terminology strong and weak.*

**Proof of the Strong Law.** First, in order to give a (relatively) short proof, we will assume

$$E[X_1^4] = L < \infty.$$

Recalling that  $E[X_1]^2 \leq E[X_1^2]$  (and so  $E[|X_1|]^2 \leq E[X_1^2]$ ), and likewise  $E[X_1^2]^2 \leq E[X_1^4]$ , we see that under this additional assumption  $E[|X_1|]$ ,  $E[X_1^2]$ , and  $E[X_1^4]$  are all finite.

Beginning with the case  $\mu = 0$ , we set

$$S_n = \sum_{j=1}^n X_j,$$

and claim that

$$E[S_n^4] = nL + 3n(n-1)E[X_1^2]^2.$$

In order to see this, we recall from the discussion of multiple combinations in Section 3.4 the notation

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

where  $n_1 + n_2 + \dots + n_N = n$ . We also recall the associated multinomial formula

$$(x_1 + x_2 + \dots + x_N)^n = \sum_{n_1 + \dots + n_N = n} \binom{n}{n_1, n_2, \dots, n_N} x_1^{n_1} x_2^{n_2} \dots x_N^{n_N}.$$

For our claim, this allows us to write

$$S_n^4 = (X_1 + X_2 + \dots + X_n)^4 = \sum_{k_1 + \dots + k_n = 4} \binom{4}{k_1, k_2, \dots, k_n} X_1^{k_1} X_2^{k_2} \dots X_n^{k_n},$$

so that (using independence)

$$E[S_n^4] = \sum_{k_1 + \dots + k_n = 4} \binom{4}{k_1, k_2, \dots, k_n} E[X_1^{k_1}] E[X_2^{k_2}] \dots E[X_n^{k_n}].$$

Now, since  $E[X_1] = 0$  (i.e.,  $\mu = 0$ ), we get no contribution from any summand with  $k_j = 1$  for some  $j \in \{1, 2, \dots, n\}$ . This means we can only get contributions from two types of terms:

- (i) Two of the  $k_j$  are 2, with the others all necessarily 0; or
- (ii) One of the  $k_j$  is 4, with the others all necessarily 0.

For (i), there are  $\binom{n}{2}$  such terms (choose the two terms that aren't 0), and for each of these the coefficient is

$$\binom{4}{2, 2, 0, \dots, 0} = \frac{4!}{2!2!} = 6,$$

with associated expectations

$$E[X_1^2] E[X_2^2] E[X_3^0] \dots E[X_n^0] = E[X_1^2] E[X_2^2] = E[X_1^2]^2.$$

(Here, the indices  $k_1$  and  $k_2$  have been distinguished for specificity, but the calculation is the same regardless of which two indices take on the values 2.) On the other hand, for (ii), there are  $\binom{n}{1} = n$  such terms, and for each of these the coefficient is

$$\binom{4}{2, 0, 0, \dots, 0} = \frac{4!}{4!} = 1,$$

with associated expectations

$$E[X_1^4] E[X_2^0] E[X_3^0] \dots E[X_n^0] = E[X_1^4].$$

In total, we see that

$$\begin{aligned} E[S_n^4] &= \binom{n}{2} 6E[X_1^2]^2 + nE[X_1^4] = \frac{n!}{2!(n-2)!} 6E[X_1^2]^2 + nL \\ &= nL + 3n(n-1)E[X_1^2]^2, \end{aligned}$$

as claimed.

We've already noted that  $E[X_1^2]^2 \leq E[X_1^4]$ , and it follows immediately that

$$E[S_n^4] \leq nL + 3n(n-1)L = L(3n^2 - 2n).$$

Dividing by  $n^4$ , we see that

$$E\left[\frac{S_n^4}{n^4}\right] \leq L\left(\frac{3}{n^2} - \frac{2}{n^3}\right).$$

The series  $\sum_{n=1}^{\infty} \frac{1}{n^2}$  and  $\sum_{n=1}^{\infty} \frac{1}{n^3}$  both converge (e.g., by the integral test), so

$$E\left[\sum_{n=1}^{\infty} \frac{S_n^4}{n^4}\right] = \sum_{n=1}^{\infty} E\left[\frac{S_n^4}{n^4}\right]$$

converges. This convergence means that the probability that the sum  $\sum_{n=1}^{\infty} \frac{S_n^4}{n^4}$  must be 1, because if not there would be some non-zero probability  $p > 0$  that it diverges, and so the expected value would include a term of the form  $p \cdot \infty = \infty$ .

We know that if a sum converges then

$$\lim_{n \rightarrow \infty} \frac{S_n^4}{n^4} = 0$$

(contrapositive to the divergence test from calculus), so this limit is a probability-1 event. I.e.,

$$P\left(\lim_{n \rightarrow \infty} \frac{S_n^4}{n^4} = 0\right) = 1,$$

which is equivalent to

$$P\left(\lim_{n \rightarrow \infty} \frac{S_n}{n} = 0\right) = 1.$$

This is the claim of the theorem in the case that  $\mu = 0$ . If  $\mu \neq 0$ , we set  $Y_j = X_j - \mu$ , so that  $E[Y_j] = 0$ , and we apply the result for  $\mu = 0$  to see that

$$P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n (X_j - \mu) = 0\right) = P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n Y_j = 0\right) = 1.$$

But  $\frac{1}{n} \sum_{j=1}^n \mu = \mu$ , so this says precisely

$$P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n X_j = \mu\right) = 1,$$

which is the assertion of the theorem. □

**Theorem 9.5.** (Central Limit Theorem) *Let  $\{X_j\}_{j=1}^{\infty}$  be a sequence of i.i.d. random variables with finite mean  $\mu = E[X_1]$  and finite standard deviation  $\sigma$ . Then the distribution of*

$$Z_n := \frac{\sum_{j=1}^n X_j - n\mu}{\sigma\sqrt{n}} = \frac{\frac{1}{n} \sum_{j=1}^n X_j - \mu}{\sigma/\sqrt{n}}$$

*approaches that of the standard normal random variable as  $n \rightarrow \infty$ . That is,*

$$\lim_{n \rightarrow \infty} P(Z_n \leq z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{x^2}{2}} dx$$

*for all  $z \in \mathbb{R}$ , with the convergence uniform in  $z$ . This means: given any  $\epsilon > 0$ , there exists  $N$  large enough so that  $n > N$  implies*

$$\left| P(Z_n \leq z) - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{x^2}{2}} dx \right| < \epsilon, \quad \forall z \in \mathbb{R}.$$

**Remark 9.3.** We won't prove the Central Limit Theorem in these notes, but let's think about what it tells us. First,

$$E[Z_n] = \frac{1}{\sigma\sqrt{n}} E\left[\sum_{j=1}^n X_j - n\mu\right] = 0$$

and

$$\text{Var}[Z_n] = \frac{1}{\sigma^2 n} \left( \text{Var}\left[\sum_{j=1}^n X_j\right] - \text{Var}[n\mu] \right) = 1.$$

Also,

$$\frac{1}{n} \sum_{j=1}^n X_j = \mu + \frac{\sigma}{\sqrt{n}} Z_n.$$

For  $n$  large,  $Z_n$  is approximately normal  $N$ , so  $\frac{1}{n} \sum_{j=1}^n X_j$  is approximately Gaussian with mean  $\mu$  and standard deviation  $\frac{\sigma}{\sqrt{n}}$ . Since  $X_j$  can have any distribution with finite mean and variance, this asserts that the Gaussian distribution is in some sense universal.

**Example 9.8.** Let's return to the trick coin from Example 7.15, and let's suppose the probability of getting heads is  $p = .6$ . As an i.i.d. sequence, we take

$$X_j = \begin{cases} 1 & \text{prob .6 (i.e., heads occurs)} \\ 0 & \text{prob .4 (i.e., tails occurs)}. \end{cases}$$

Then  $\mu = E[X_j] = .6$ , and we also have  $E[X_j^2] = .6$  and

$$\sigma^2 = E[X_j^2] - E[X_j]^2 = .6 - .36 = .24.$$

The Central Limit Theorem asserts that

$$Z_n = \frac{\sum_{j=1}^n X_j - n\mu}{\sigma\sqrt{n}} = \frac{\sum_{j=1}^n X_j - .6n}{\sqrt{.24n}}$$

will be approximately standard normal as  $n$  becomes large.

In order to check this, we'll use the MATLAB M-file *central1.m*, included below.

```
function central1(p,n,N)
%CENTRAL1: Simulates a binomial random variable
%given its probability p and the number of
%trials N, then compares the results with
%a standard normal distribution.
%p = probability of heads
%n = number of flips for each simulation
%N = number of simulations of Z_n
%
rng('shuffle')
mu = p; v = p-p^2; %expected value and variance
```



```

for k=1:N
X = round(rand([1,n])+(p-1/2));
Z(k) = (sum(X)-n*mu)/(sqrt(v*n));
end
%
%Comparison with standard normal
[counts,centers]=hist(Z,50); %Histogram with 50 bins
hist(Z,50)
hold on
pause
scale = (centers(2)-centers(1))*N; %Area of histogram = binwidth*Total num-
ber of points
x = linspace(-4,4,1000);
f=scale/sqrt(2*pi)*exp(-x.^2/2);
plot(x,f,'r','LineWidth',2)

```

If we implement this file with the command `>> central1(.6,5000,10000)`, corresponding with taking  $p = .6$ ,  $n = 5000$ , and simulating  $Z_{5000}$  10000 times, we obtain the histogram in Figure 9.3, overlaid with the expected Gaussian distribution.  $\triangle$

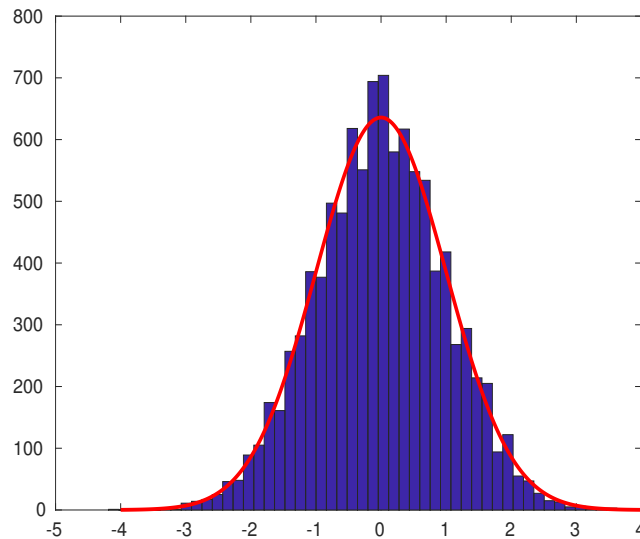


Figure 9.3: Histogram and (scaled) Gaussian distribution for Example 7.16.

We often apply the Central Limit Theorem as follows: suppose we want to estimate a probability  $p$  with the estimator

$$\hat{p} := \frac{1}{n} \sum_{j=1}^n X_j.$$

This could be a simulation problem such as the ones discussed above, or it could be a data collection problem along the lines of those discussed in Section 8 below. As an example of

such problems, we might want to estimate the probability that a certain person has a disease, given some test result. In this case, our i.i.d. random variables would naturally have the Bernoulli form

$$X_j = \begin{cases} 1 & \text{if person } j \text{ has the disease} \\ 0 & \text{if person } j \text{ does not have the disease.} \end{cases}$$

If  $n$  is large, the Central Limit Theorem asserts that the random variable

$$Z_n = \frac{\hat{p} - \mu}{\sigma/\sqrt{n}}$$

is approximately standard normal  $N$ , where in the Bernoulli case

$$\begin{aligned} \mu &= E[X_j] = p \\ \sigma^2 &= \text{Var}[X_j] = p - p^2. \end{aligned}$$

So, for any  $k > 0$ ,

$$P\left(\left|\frac{\hat{p} - p}{\sqrt{(p - p^2)/n}}\right| \leq k\right) \cong \frac{1}{\sqrt{2\pi}} \int_{-k}^{+k} e^{-\frac{x^2}{2}} dx.$$

In particular, this integral gives the probability that

$$p - k\sqrt{\frac{p - p^2}{n}} \leq \hat{p} \leq p + k\sqrt{\frac{p - p^2}{n}}.$$

Although,  $k$  can be any positive value, integer values are most common, especially:

$$\begin{aligned} k = 1 &: \frac{1}{\sqrt{2\pi}} \int_{-1}^{+1} e^{-\frac{x^2}{2}} dx = .6827 \\ k = 2 &: \frac{1}{\sqrt{2\pi}} \int_{-2}^{+2} e^{-\frac{x^2}{2}} dx = .9545. \end{aligned}$$

We are about 68% certain that

$$\hat{p} \in \left[p - \sqrt{\frac{p - p^2}{n}}, p + \sqrt{\frac{p - p^2}{n}}\right],$$

and about 95% certain that

$$\hat{p} \in \left[p - 2\sqrt{\frac{p - p^2}{n}}, p + 2\sqrt{\frac{p - p^2}{n}}\right].$$

Here,  $\hat{p}$  is a random variable, and  $p$  is unknown, so what does this give us? In practice, we often take the experimentally measured value  $\hat{p}^*$  as an approximation of  $p$ , and use confidence intervals of the form

$$\left[\hat{p}^* - k\sqrt{\frac{\hat{p}^* - (\hat{p}^*)^2}{n}}, \hat{p}^* + k\sqrt{\frac{\hat{p}^* - (\hat{p}^*)^2}{n}}\right].$$

We take this to mean that the value of the parameter  $p$  that we are trying to determine is

$$p = \hat{p}^* \pm \sqrt{\frac{\hat{p}^* - (\hat{p}^*)^2}{n}}.$$

**Example 9.9.** For the coin described in Examples 9.7 and 9.8, determine the number of simulations required to ensure with 95% confidence (i.e., with  $k = 2$ ) that the experimentally determined value  $\hat{p}^*$  is within a maximum error of .01 from  $p$ ; i.e., that  $|p - \hat{p}^*| \leq .01$ .

For this, we require

$$2\sqrt{\frac{\hat{p}^* - (\hat{p}^*)^2}{n}} \leq .01,$$

and since  $\hat{p}^* - (\hat{p}^*)^2 \leq \frac{1}{4}$ , this is

$$\frac{1}{\sqrt{n}} \leq .01 \implies \sqrt{n} \geq \frac{1}{.01} \implies n \geq \frac{1}{.0001} = 10,000.$$

I.e., if we use  $n = 10,000$  and compute  $\hat{p}^*$ , we will have a 95% probability that our error will be at most .01. △

**Example 9.10.** For the coin described in Examples 9.7, 9.8, and 9.9, simulate  $n = 10,000$  flips, and compute the resulting value  $\hat{p}^*$ , along with 68% and 95% confidence intervals.

We carry out the simulations with MATLAB, using *bern.m*.

```
function bern(p,n)
%BERN: Simulates a Bernoulli random variable
%given its probability p and the number of
%trials n, then uses the simulation to
%estimate p and determine confidence
%intervals.
%
rng('shuffle')
X = round(rand([1,n])+(p-1/2));
%Estimator
phat = sum(X)/n
%
exacterror = abs(phat-p);
k1conf = sqrt((phat-phat^2)/n);
k2conf = 2*sqrt((phat-phat^2)/n);
fprintf('The exact error is %5.4f\n',exacterror)
fprintf('With 68%% confidence, the error is at most %5.4f\n',k1conf)
fprintf('With 95%% confidence, the error is at most %5.4f\n',k2conf)
```

The output from a particular run is as follows:

```
bern(.6,10000)
phat =
0.5946
The exact error is 0.0054
With 68% confidence, the error is at most 0.0049
With 95% confidence, the error is at most 0.0098
```

We see that  $\hat{p}^* = .5964$ . With 68% confidence, we can report

$$p = .5964 \pm .0049,$$

while with 95% confidence we can report

$$p = .5964 \pm .0098.$$

In this case, the actual error falls outside the 68% confidence interval, but within the 95% confidence interval.

## 10 Hypothesis Testing

Once we have determined which probability density function appears to best fit our data, we need a method for testing how good the fit is. In general, the analysis in which we test the validity of a certain statistic is called *hypothesis testing*. Before considering the case of testing an entire distribution, we will work through a straightforward example involving the test of a mean value.

### 10.1 General Hypothesis Testing

**Example 10.1, Part 1.** In the spring of 2003 the pharmaceutical company VaxGen published the results of a three-year study on their HIV vaccination. The study involved 5,009 volunteers from the United States, Canada, Puerto Rico, and the Netherlands. Overall, 97 out of 1679 placebo recipients became infected, while 190 out of 3330 vaccine recipients became infected. Of the 498 non-white, non-hispanic participants, 17 out of 171 placebo recipients became infected while 12 out of 327 vaccine recipients became infected. Determine whether it is reasonable for VaxGen to claim that their vaccination is successful.

Let  $N$  denote the number of participants in this study who were vaccinated ( $N = 3330$ ), and let  $p_0$  denote the probability of a placebo recipient becoming infected ( $p_0 = \frac{97}{1679} = .058$ ). (Note that  $p_0$  is the probability of infection in the absence of vaccination. We expect the probability  $p$  of a vaccine recipient becoming infected to be less than  $p_0$ .) Next, consider the possibility of repeating exactly the same study with a different set of 3330 participants, and set

$$X_j = \begin{cases} 1 & \text{if participant } j \text{ becomes infected} \\ 0 & \text{if participant } j \text{ does not become infected.} \end{cases}$$

This puts us in the setting of Examples 9.7 through 9.10 from the previous section, and we can proceed similarly as there. For this, we denote by  $\hat{p}$  the random variable

$$\hat{p} = \frac{1}{N} \sum_{j=1}^{3330} X_j.$$

The goal of hypothesis testing in this case is to determine how representative  $p_0$  is of the values that the random variable  $\hat{p}$  can assume, and so in particular of the actual (unknown)

value  $p$  that a vaccinated participant will become infected. (VaxGen would like to demonstrate that  $p < p_0$ , hence that the vaccine is effective.) Our benchmark hypothesis, typically referred to as the *null hypothesis* and denoted  $H_0$ , is that the vaccine is *not* effective. That is,

$$H_0 : \quad p = p_0 = .058.$$

We test our null hypothesis against our *alternative hypothesis*, typically denoted  $H_1$ , that the vaccine *is* effective. That is,

$$H_1 : \quad p < p_0.$$

Here's the main idea. We are going to assume  $H_0$ ; i.e., that  $p = p_0 = .058$  is fixed. Our experimental observation, however, is that  $\hat{p}^* = \frac{190}{3330} = .057$ , a little better. We will determine the probability that our random sample determined  $\hat{p} \leq .057$  given that the true underlying value of  $p$  is  $p_0$ . If this is highly unlikely, we reject the null hypothesis.<sup>18</sup>

Using the Central Limit Theorem, we saw in the previous section that the random variable

$$Z_N := \frac{\hat{p} - p}{\sqrt{(p - p^2)/N}}$$

approximately satisfies a standard normal distribution. Re-writing this as

$$\hat{p} = p + \sqrt{\frac{p - p^2}{N}} Z_N,$$

we see that  $\hat{p}$  is a Gaussian random variable with mean  $\mu = p$  and standard deviation  $\sigma = \sqrt{(p - p^2)/N}$ .

If  $p = p_0$  (the null hypothesis), then

$$\mu = p_0 = .058, \quad \text{and} \quad \sigma = \sqrt{\frac{p - p^2}{N}} = \sqrt{\frac{.058(1 - .058)}{3330}} = .004.$$

We can compute probabilities on  $\hat{p}$  from the Gaussian distribution

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

In particular, we are interested in the probability that  $\hat{p} \leq .057$ . We have,

$$P(\hat{p} \leq .057) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{.057} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = .4013,$$

where the integration has been carried out with MATLAB. We see that we have a 40% chance of getting  $\hat{p} \leq .057$  even if the vaccine is not effective at all. This is *not* good enough to reject  $H_0$  and we conclude that VaxGen cannot legitimately claim that their vaccine is effective.

---

<sup>18</sup>Think about flipping a fair coin ten times and counting the number of times it lands heads. If the coin is fair,  $p_0 = .5$ , but for each ten-flip experiment  $p_1 = \frac{\text{Number of heads}}{10}$  will be different. If  $p_1 = .1$ , we might question whether or not the coin is genuinely fair.

**Example 10.1, Part 2.** When this study came out, VaxGen understandably played down the main study and stressed the results in the non-white, non-Hispanic subset of participants, for which  $p_0 = \frac{17}{171} = .099$  and  $\hat{p}^* = \frac{12}{327} = .037$ . In this case,

$$\mu = p_0 = .099 \quad \text{and} \quad \sigma = \sqrt{\frac{.099(1 - .099)}{327}} = .0165,$$

from which we compute

$$P(\hat{p} \leq .0370) = .00009.$$

The finding, then, is that the probability that  $\hat{p} \leq .037$  is .009% and we reject the null hypothesis and claim that the vaccination is effective in this subset of participants.<sup>19</sup>  $\triangle$

A critical question becomes, how low does the probability of the event have to be before the null hypothesis is rejected. In Example 8.1, Part 1, the probability was over 40%, not that much better than a coin toss. We certainly cannot reject it based on that. In Example 8.1, Part 2, the probability was .009%: an anomaly like this might occur 9 times in 100,000 trials. In this case, we are clearly justified in rejecting the null hypothesis. What about cases in between: 10%, 5%, 1%? In the end, decisions like this are largely made on the requirement of accuracy.

## 10.2 Hypothesis Testing for Distributions

In the case of hypothesis testing for distributions, our idea will be to check our theoretical cumulative distribution function against the *empirical distribution function*.

### 10.2.1 Empirical Distribution Functions

For a random variable  $X$  and a set of  $n$  observations  $\{x_1, x_2, \dots, x_n\}$ , we define the *empirical distribution function*  $F_e(x)$  by

$$F_e(x) = \frac{\text{Number of } x_k \leq x}{n}.$$

**Example 10.2.** Consider again the Lights, Inc. data given in Table 7.1, and let  $F_e(x)$  denote the empirical distribution function for the times to failure. Zero lights failed between 0 and 400 hours, so

$$F_e(400) = \frac{0}{n} = 0.$$

By 500 hours 2 lights had failed, while by 600 hours 5 lights had failed, and we have

$$F_e(500) = \frac{2}{100} = .02$$

$$F_e(600) = \frac{5}{100} = .05.$$

The MATLAB M-file *edf.m* takes a vector of observations and a point and calculates the empirical distribution of those observations at that point.

---

<sup>19</sup>Though several factors suggest that significantly more study is necessary, since the sample size is so small in this subset, and since in any set of data, subsets will exist in which results are favorable.

```

function f = edf(x,D);
%EDF: Function file which returns the empirical
%distribution function at a point x of a set
%of data D.
f = sum(D <= x)/length(D);

```

If  $T$  denotes the data for Lights, Inc., the following MATLAB diary file shows the usage for *edf.m*. (Recall that the data for Lights, Inc. is recorded in the MATLAB M-file *lights.m*.)

```

>>lights
>>edf(400,T)
ans =
0
>>edf(500,T)
ans =
0.0200
>>edf(600,T)
ans =
0.0500
>>edf(1500,T)
ans =
0.9800
>>edf(1600,T)
ans =
1

```

Our theoretical distribution for this data was Gaussian with  $\mu = 956.88$  and  $\sigma = 234.69$ , corresponding with the cumulative distribution function

$$F(x) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^x e^{-\frac{(y-\mu)^2}{2\sigma^2}} dy.$$

The MATLAB script M-file *edfplot.m* compares the theoretical distribution with the empirical distribution (see Figure 10.1).

```

%EDFPLOT: MATLAB script M-file for plotting
%the EDF along with the CDF for Lights, Inc. data.
x=linspace(400,1600,100);
for k=1:length(x);
Fe(k)=edf(x(k),T);
F(k)=quad('1/(sqrt(2*pi)*234.69)*exp(-(y-956.88).^2/(2*234.69^2))',0,x(k));
end
plot(x,Fe,x,F,'-')

```

Though certainly not perfect, the fit of our Gaussian distribution at least seems plausible. Determining whether or not we accept this fit is the subject of *hypothesis testing*.  $\triangle$

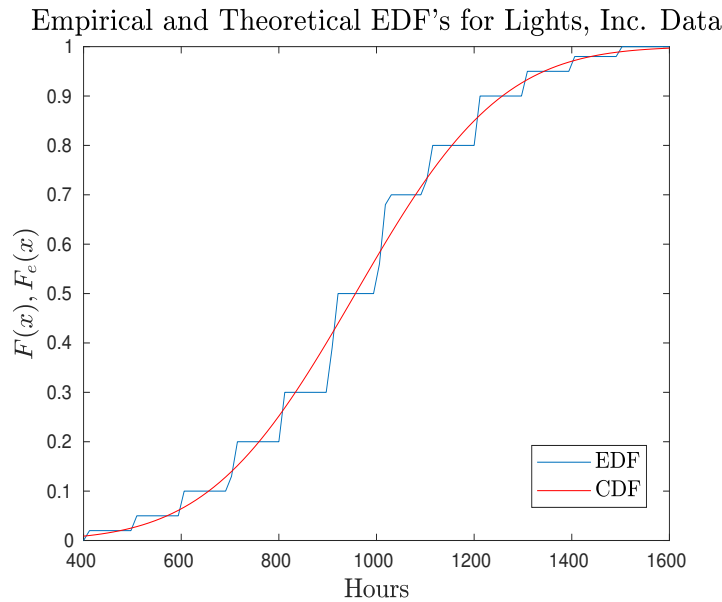


Figure 10.1: Empirical and Theoretical cumulative distribution functions for Lights, Inc. data.

**Example 10.3.** Consider again our data from Lights, Inc., summarized in Table 7.1 and denoted by the vector  $T$ . After looking at a histogram of this data, we determined that it was well described by a Gaussian distribution with  $\mu = 956.88$  and  $\sigma = 234.69$ . Here, we consider whether or not this distribution is indeed a reasonable fit.

First, we require some quantifiable measure of how closely our distribution fits the data. (So far we've simply been glancing at the data and judging by the shape of its histogram.) One method for doing this is to compare the proposed cumulative distribution function,  $F(x)$ , with the data's empirical distribution function,  $F_e(x)$  (see Figure 10.1). Two standard tests are the *Kolmogorov–Smirnov* test and the *Cramer–von Mises* test.

**1. Kolmogorov–Smirnov statistic.** The simplest test statistic for cumulative distribution functions is the *Kolmogorov–Smirnov* statistic, defined by

$$D := \sup_{x \in \mathbb{R}} |F_e(x) - F(x)|.$$

Referring, for example, to Figure 10.1, the Kolmogorov–Smirnov statistic simply follows the entire plot of both functions and determines the greatest distance between them. In the event that  $F_e$  and  $F$  have been defined as vectors, we can compute the K–S statistic in MATLAB with the single command  $D = \max(\text{abs}(F_e - F))$ . For our Lights, Inc. data  $D = .0992$ .<sup>20</sup>

**2. Cramer–von Mises statistic.** A second test statistic that measures the distance between  $F(x)$  and  $F_e(x)$  along the entire course of the functions is called the *Cramer–von*

<sup>20</sup>Observe that  $D$  can be calculated more accurately by refining  $x$ . We can accomplish this by taking more points in the *linspace* command.



*Mises* statistic and is given by

$$W^2 := N \int_{-\infty}^{+\infty} (F_e(x) - F(x))^2 f(x) dx,$$

where  $N$  is the number of data points and  $f(x) = F'(x)$  is the (proposed) probability density function. The C–vM statistic is slightly more difficult to analyze than the K–S statistic, largely because the empirical distribution function,  $F_e(x)$ , can be cumbersome to integrate. The primary thing to keep in mind is that integrands in MATLAB must accept vector data. We compute the integrand for the C–vM test with the MATLAB M-file *cvm.m*. Observe that the command *lights* simply defines the vector  $T$  (i.e., we could have replaced this command with  $T=[415,478,\dots]$ , the entire vector of data).

```
function value = cvm(x,T)
%CVM: Returns the Cramer-von Mises integrand
%for the given data and distribution.
mu = 956.88; sig = 234.69;
f = @(y) 1/(sqrt(2*pi)*sig)*exp(-(y-mu).^2/(2*sig^2));
for k=1:length(x)
F(k)=integral(f,0,x(k));
Fe(k) = edf(x(k),T);
fval(k) = f(x(k));
end
value = (F-Fe).^2.*fval;
```

We now compute  $W^2$  as  $Wsq=length(T)*integral(@(x) cvm(x,T),400,1600)$ . We find  $W^2 = .1370$ .

Finally, we use the K–S and C–vM statistics to test the adequacy of our proposed distribution. Our null hypothesis in this case is that our proposed distribution adequately describes the data,

$$H_0 : F_d(x) = F(x),$$

while our alternative hypothesis is that it does not,

$$H_1 : F_d(x) \neq F(x).$$

We test our null hypothesis by testing one of the test statistics described above. In order to accomplish this, we first need to determine the (approximate) distribution of our test statistic. (In our HIV example, our test statistic was  $p$  and its approximate distribution was Gaussian.)

Let's focus first on the K–S statistic,  $D$ . Observe that  $D$  is a random variable in the following sense: each time we test 100 bulbs, we will get a different outcome and consequently a different value for  $D$ . Rather than actually testing more bulbs, we can simulate such a test, assuming  $H_0$ , i.e., that the data really is arising from a Gaussian distribution with  $\mu = 956.88$  and  $\sigma = 234.69$ . In the MATLAB M-file *kstest.m* we simulate a set of data points and compute a new value of  $D$ . Observe that the vectors  $x$  and  $F$  have already been computed as above, and  $T$  contains the original data.

```

function value = kstest(F,x,T)
%KSTEST: Takes a vector cdf F and a vector
%x and original data T and determines
%the Kolmogorov-Smirnov
%statistic for a sample taken from a
%normal distribution.
clear G; clear Fe;
N = length(F); %Number of points to consider
%Simulate Gaussian data
mu = mean(T);
sd = std(T);
for k=1:length(T) %Always simulate the same number of data points in ex-
periment
    G(k)=mu+randn*sd;
end
%Compute empirical distribution function
for k=1:N
    Fe(k)=edf(x(k),G);
end
%Kolmogorov-Smirnov statistic
value = max(abs(Fe-F));

```

Working at the MATLAB Command Window prompt, we have

```

kstest(F,x,T)
ans =
0.0492
kstest(F,x,T)
ans =
0.0972
kstest(F,x,T)
ans =
0.0655

```

We next develop a vector of  $D$  values and consider a histogram (see Figure 10.2).

```

>>for j=1:1000
D(j)=kstest(F,x,T);
end
>>hist(D,50)

```

Of the distributions we've considered, this data is best fit by either a beta distribution or a Weibull distribution. One method by which we can proceed is to develop an appropriate PDF for this new random variable  $D$ , and use it to compute the probabilities required for our hypothesis testing. In fact, this  $D$  is precisely the random variable  $D$  we analyzed in our section on maximum likelihood estimation. Observing, however, that we can create as many

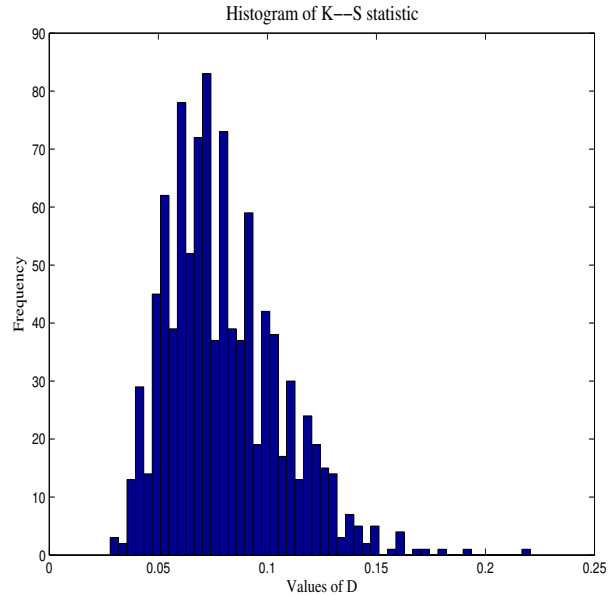


Figure 10.2: Histogram for K–S statistic, Lights, Inc. data.

realizations of  $D$  as we like, and that as the number of observations increases, the empirical distribution function approaches the continuous distribution function, we proceed directly from the EDF. That is, we compute

$$P(D \geq .0992) \cong 1 - F_e(.0992).$$

In MATLAB, the command `1-edf(.0992,D)` returns .23, or 23%. This signifies that if our distribution is correct, there remains a 23% chance that our statistic  $D$  will be as large as it was. Though we would certainly like better odds, this is typically considered acceptable.  $\triangle$

## 11 Application to Finance

In this section, we introduce some of the basic mathematical tools involved in the study of finance.

### 11.1 Random Walks

Suppose two individuals flip a fair coin each day, and exchange money according to the following rule: if the coin lands with heads up, player A pays player B one dollar, while if the coin lands tails up, player B pays player A one dollar. We define the series of random variables

$$X_t = \text{Player A's earnings on day } t.$$

For example, we have

$$X_1 = X = \begin{cases} +1 & \text{with probability } 1/2 \\ -1 & \text{with probability } 1/2 \end{cases},$$

while

$$X_2 = X_1 + X = \begin{cases} +2 & \text{with probability } 1/4 \\ 0 & \text{with probability } 1/2 \\ -2 & \text{with probability } 1/4 \end{cases},$$

where  $X$  denotes the random event carried out each day and the pattern continues with  $X_{t+1} = X_t + X$ . Defined as such,  $X_t$  is a *random walk*; more generally, any process such as  $X_t$  for which each value of  $t$  corresponds with a different random variable will be referred to as a *stochastic process*. We make the following critical observations regarding  $X_t$  from the simple game described above:

1.  $X_{t+1} - X_t$  is independent of  $X_t - X_{t-1}$ , which in turn is independent of  $X_{t-1} - X_{t-2}$  etc. That is, the coin toss on any given day is entirely independent of all preceding coin tosses.
2.  $E[X_t] = 0$  for all  $t$ .
3.  $\text{Var}[X_t] = t$  for all  $t$ .

Also, we observe that the recursive nature of  $X_t$  makes it particularly easy to simulate with MATLAB. For example, the following M-file serves to simulate and plot a random walk associated with this game.

```
%RANWALK: MATLAB M-file written for the purpose
%of simulating and plotting a random walk.
clear X; %Initialize random variable
N = 100; %Number of steps in the random walk
X(1) = 0; %Start at 0
for k=1:N
X(k+1) = X(k) + 2*randi(2)-3;
end
t = 0:N;
plot(t,X)
title('Random Walk Example','interpreter','latex','FontSize',16)
xlabel('Number of flips','interpreter','latex','FontSize',14)
ylabel('Location of the Random Walk','interpreter','latex','FontSize',14)
```

Two plots generated with `ranwalk.m` are shown in Figure 11.1.

## 11.2 Brownian Motion

In 1905, at the age of 26, a little known Swiss patent clerk in Berne published four landmark papers: one on the photoelectric effect (for which he would receive a Nobel prize in 1921), two on the theory of special relativity, and one on the transition density for a phenomenon that had come to be known as Brownian motion. The clerk was of course Albert Einstein, so we're in good company with our studies here.<sup>21</sup>

---

<sup>21</sup>Though Einstein was working in Switzerland at the time, he was born in Ulm, Germany.

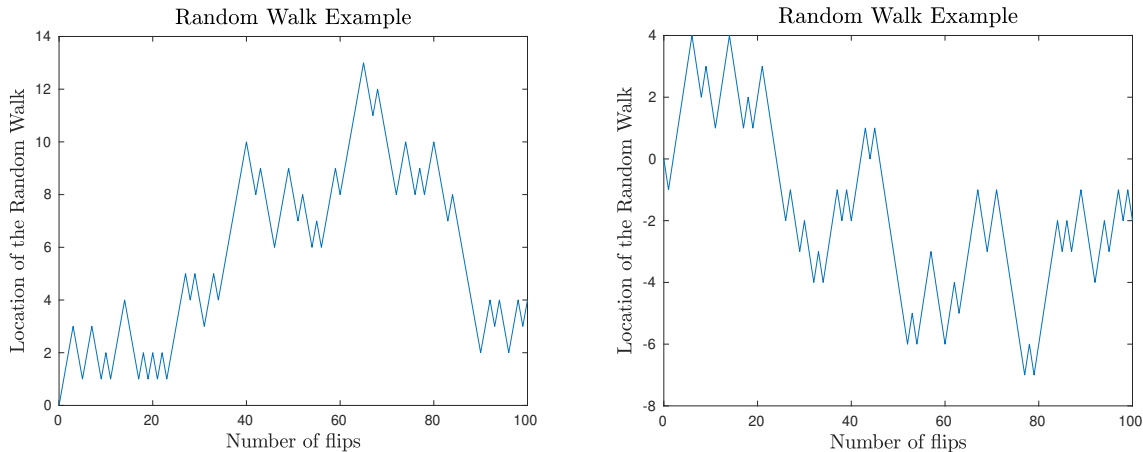


Figure 11.1: Random walk over 100 days.

Brownian motion is the name given to the irregular movement of pollen, suspended in water, observed by the Scottish botanist Robert Brown in 1828. As Einstein set forth in his paper of 1905, this movement is caused by collisions with water molecules. The dynamics are so complicated that the motion appears completely random. Our eventual goal will be to use the framework of mathematics built up around Brownian motion to model “random” behavior in the stock market.

A good intuitive way to think of Brownian motion is as a continuous time random walk. The random walk discussed in Section 8.1 is carried out at discrete times, one flip each day. Of course, we could increase the number of flips to say 2 each day, or 4 or 8. As the number of flips increases, the process looks more and more like a continuous process: at some point we find ourselves doing nothing else besides flipping this stupid coin. Rather, however, than actually flip all these coins, we’ll simply build Brownian motion from a continuous distribution, the normal distribution. As it turns out, these two approaches—infinite flipping and normal distributions—give us exactly the same process (though we will make no effort to prove this).

The particular definition we’ll use is due to Norbert Wiener (1894–1964), who proposed it in 1918. (Compare with Properties 1, 2, and 3 from Section 8.1.)

**Definition 11.1.** (Standard Brownian Motion)<sup>22</sup> *A stochastic process  $B_t$ ,  $t \geq 0$ , is said to be a standard Brownian motion if:*

1.  $B_0 = 0$ .
2.  $B_{t_n} - B_{t_{n-1}}, B_{t_{n-1}} - B_{t_{n-2}}, \dots, B_{t_2} - B_{t_1}, B_{t_1}$  are all independent with distributions that depend (respectively) only on  $t_n - t_{n-1}, t_{n-1} - t_{n-2}, \dots, t_2 - t_1, t_1$ ; that is, only on the time interval between observations.
3. For each  $t \geq 0$ ,  $B_t$  is normally distributed with mean 0 and variance  $t$ . ( $E[B_t] = 0$ ,  $\text{Var}[B_t] = t$ .)

---

<sup>22</sup>We omit entirely any attempt to prove that a process satisfying this definition exists. See, e.g., Chapter 2 in [KS1991] for details.

Critically, Property 3 is equivalent to the assertion that the probability density function for a standard Brownian motion  $B_t$  is

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

Hence, we can carry out all the usual analyses. For example, to compute the probability that  $B_t$  lies on the interval  $[a, b]$ , we compute

$$P(a \leq B_t \leq b) = \frac{1}{\sqrt{2\pi t}} \int_a^b e^{-\frac{x^2}{2t}} dx.$$

Just as with random walks, Brownian motion is straightforward to simulate using MATLAB. Recalling from Section 9.3 that MATLAB's command `randn` creates a normally distributed random variable with mean 0 and variance 1, we observe that  $\sqrt{t} * \text{randn}$  will create a random variable with mean 0 and variance  $t$ . That is, if  $N$  denotes a standard normal random variable (the random variable simulated by `randn`), we have

$$\begin{aligned} E[\sqrt{t}N] &= \sqrt{t}E[N] = 0 \\ \text{Var}[\sqrt{t}N] &= E[tN^2] - E[\sqrt{t}N]^2 = tE[N^2] = t. \end{aligned}$$

Our fundamental building block for Brownian paths will be the MATLAB function M-file `brown.m`, which will take a time  $t$  and return a random value associated with that time from the appropriate distribution.

```
function value = brown(t);
%BROWN: Returns the value of a standard Brownian motion
%at time t.
value = sqrt(t)*randn;
```

A *Brownian path* is simply a series of realizations of the Brownian motion at various times  $t$ . For example, the following M-file `snbm.m` describes a Brownian path over a period of 50 days, with a realization at each day.

```
%SNBM: Simulates a standard normal Brownian motion.
T = 100; %Number of time steps
delT = 1; %Time increment
time = 0; %Starting time
t(1) = 0; %Begin at t = 0
B(1) = 0; %Corresponds with t = 0
for k = 2:T+1;
t(k) = t(k-1) + delT;
B(k) = B(k-1) + brown(delT);
end
plot(t,B)
title('Example Brownian Path','interpreter','latex','FontSize',16)
xlabel('Number of Time Steps','interpreter','latex','FontSize',14)
ylabel('Location of the Brownian Path','interpreter','latex','FontSize',14)
```

Two Brownian paths generated with `snbm.m` are shown in Figure 11.2.

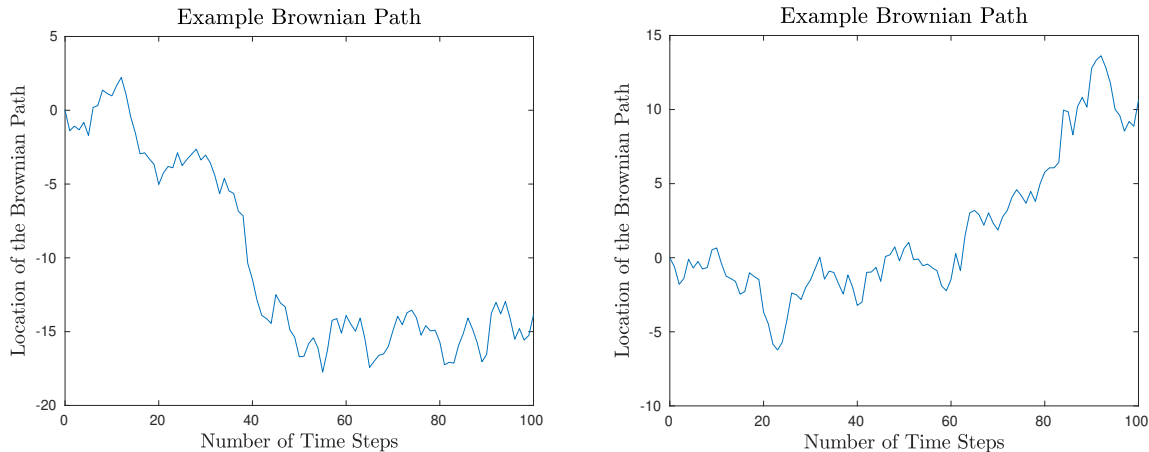


Figure 11.2: Standard Brownian motion over 100 days.

### 11.3 Stochastic Differential Equations

In the event that there was nothing whatsoever deterministic about the stock market, Brownian motion alone would suffice to model it. We would simply say that each day the probability of the stock's value rising a certain amount is exactly the same as the probability of its declining by the same amount. If this were the case, however, the stock market would most likely have closed a long time ago.

A popular and reasonably effective measure of overall stock market behavior is the Dow Jones industrial average, initiated by the journalist Charles Dow in 1884, along with his fellow journalist and business partner Edward Jones. In 1882, Dow and Jones formed Dow Jones & Company, which gathered and sold news important to financial institutions (this newsletter would later morph into the Wall Street Journal.) The stock market was relatively new and untested at the time, and one thing companies wanted was an idea of how it was behaving in general. Dow's idea was to follow 11 important stocks (mostly railroads back then) in hopes of finding overall market trends. On October 1, 1928, long after Dow's death, the number of stocks in his average was fixed at 30, the current number. Along with a number of companies now defunct, this list included such familiar names as Chrysler, General Electric, and Sears, Roebuck. At the end of 1928, the DJIA (Dow Jones Industrial Average) stood at roughly 300. Today (March 11, 2022), it sits at 33,189.89. Over 94 years this corresponds with an average yearly increase of 5.13%.<sup>23</sup>

In order to introduce this kind of deterministic growth into our model for a stock, we might simply add a term that corresponds with this growth. Recalling that continuously compounded interest grows like  $Pe^{rt}$ , where  $P$  is the principal investment (cf. Footnote 23) we have the stock price model

---

<sup>23</sup>The reader may recall from high school algebra that a principal investment  $P$ , invested over  $t$  years at interest rate  $r$  yields a return  $R_t = P(1+r)^t$ . Hence, the equation I've solved here is  $33189.89 = 300(1+r)^{94}$ . Typically, the numbers you'll hear bandied around are more in the ballpark of 10% or 11%; for example, in his national bestseller, *A Random Walk Down Wall Street*, Burton G. Malkiel points out that between 1926 and 1994, the average growth on all common stocks has been 10.2%. The point here is simply that, historically speaking, stock behavior is not entirely random: these things are going up.

$$S_t = S_0 e^{rt} + B_t,$$

determined growth corrected by a random fluctuation.

The next step requires some motivation. In many cases, it's easier to express an equation for the rate of change of some quantity than it is to write down an expression for the quantity itself.<sup>24</sup> This suggests that it might be useful to introduce an analogue to differential equations in the context of stochastic processes. Though a quick glance at the jagged turns in Figures 11.1 and 11.2 might lead us (correctly) to suspect that stochastic processes don't generally have classical derivatives, we can (for the sake of all that \$\$ to be made gambling on the stock market) define the *differential form* (i.e., in this context, a multidimensional polynomial of differentials)

$$dS_t = S_0 r e^{rt} dt + dB_t, \quad (11.1)$$

where we view  $dt$  as a small but finite increment in time. Equations of form (11.1) are called *stochastic differential equations*.

More generally, the stochastic differential equation for a reasonably well-behaved stochastic process  $X_t$  can be written in the form

$$dX_t = \mu(t, X_t)dt + \sigma(t, X_t)dB_t. \quad (11.2)$$

Though such equations can be analyzed exactly, we will focus on solving them numerically. We recall from introductory courses in differential equations that Euler's numerical method for solving a first order equation  $x'(t) = f(t, x)$  is derived by approximating  $x'(t)$  with a difference quotient:

$$x'(t) = \lim_{h \rightarrow 0} \frac{x(t+h) - x(t)}{h} \Rightarrow x'(t) \cong \frac{x(t+h) - x(t)}{h}, \quad h \text{ small.}$$

Euler's approximation becomes, then,

$$\frac{x(t+h) - x(t)}{h} = f(t, x) \Rightarrow x(t+h) = x(t) + hf(t, x),$$

an iterative equation ideal for computation. For (11.2) we note that  $dt$  plays the role of  $h$  and  $dB_t$  is the critical new issue. We have

$$X_{t+dt} = X_t + \mu(t, X_t)dt + \sigma(t, X_t)(B_{t+dt} - B_t).$$

Recalling that  $B_{t+dt} - B_t \stackrel{d}{=} B_{dt}$ , ( $\stackrel{d}{=}$  means equal in probability, that the quantities on either side of the expression have the same probability density function), our algorithm for computation will be

$$X_{t+dt} = X_t + \mu(t, X_t)dt + \sigma(t, X_t)B_{dt},$$

where  $B_{dt}$  will be recomputed at each iteration. Assuming appropriate MATLAB function M-files *mu.m* and *sig.m* have been created (for  $\mu$  and  $\sigma$  respectively), the following function M-file will take a time interval  $t$  and initial value  $z$  and solve (11.2) numerically.

---

<sup>24</sup>Hence the field of differential equations.



```

function sdeplot(t,z)
%SDEPLOT: Employs Euler's method to numerically
%an SDE, and to plot the solution. The drift is
%computed in the subfunction mu(), and the diffusion
%is computed in the subfunction sig().
%The variable t is time, and z is initial position.
clear time;
clear X;
steps = 1000; %Number of increments
dt = t/steps; %time increment
X(1) = z; %Set initial data
m = 1; %Initialization of vector entry
time(1) = 0;
for k=dt:dt:t
X(m+1) = X(m) + mu(k,X(m))*dt + sig(k,X(m))*sqrt(dt)*randn;
time(m+1) = time(m) + dt;
m = m + 1;
end
plot(time,X)
title('Example SDE solution','interpreter','latex','FontSize',16)
xlabel('Time $t$','interpreter','latex','FontSize',14)
ylabel('$X_t$','interpreter','latex','FontSize',14)
%
function f = mu(t,x);
%MU: Contains drift term for a stochastic
%differential equation.
r = .0513;
f = r*x;
function f = sig(t,x);
%SIG: Contains diffusion term for a stochastic
%differential equation.
s=.05;
f = s*x;

```

Two solution paths computed with the input `sdeplot(94,300)`, and the SDE

$$dX_t = .0513X_t dt + .05X_t dB_t$$

are shown in Figure 11.3. These roughly follow the DJIA during the years 1928 to 2022.

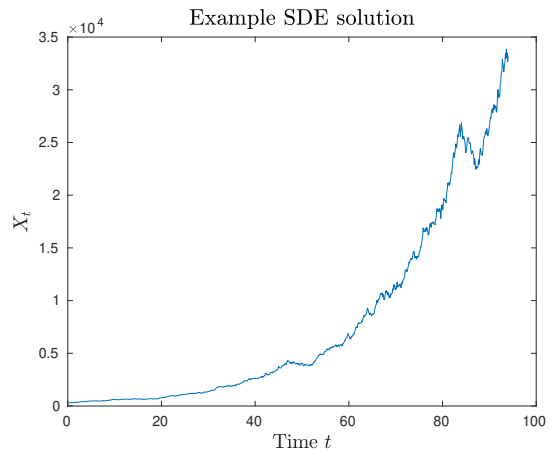
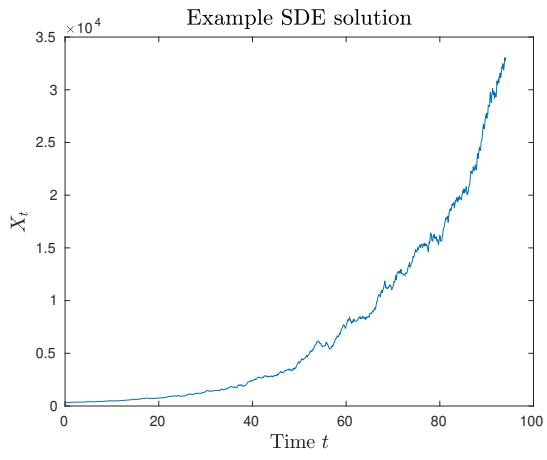


Figure 11.3: Geometric Brownian Motion

# Index

- alternative hypothesis, 93
- Bernoulli probability mass function, 33
- beta distribution, 60
- beta function, 60
- binomial probability mass function, 33
- bins, 55
- brownian motion, 100
  
- complement, 4
- Conditional probability, 15
- containment, 4
- Covariance, 30
- Cramer–von Mises statistic, 96
- cumulative distribution function, 31
  
- dominance principle, 36
- Dow Jones, 103
  
- empirical distribution function, 94
- Equations, 103
- event, 3
- expected value, 23
- exponential distribution, 58
  
- game theory, 34
- gamma distribution, 61
- gamma function, 60
- Gaussian distribution, 56
- geometric probability mass function, 33
  
- intersection, 3
  
- Kolmogorov–Smirnov statistic, 96
  
- likelihood function, 70
  
- MATLAB commands
  - beta(), 61
  - gamma(), 60
  - hist(), 55
  - rand, 74
- MATLAB functions
  - randn, 76
- minimax method, 40
  
- mixture distributions, 61
- Monte Carlo method, 73
  
- normal distribution, 56
- null hypothesis, 93
  
- oddmens, 45
- outcome, 3
  
- Poisson probability mass function, 33
- probability mass function, 32
- pseudo random variables, 74
  
- queueing theory, 80
  
- randi(), 76
- random variables
  - continuous, 51
  - discrete, 22
- random walk, 100
- realizations, 22
- rejection method, 77
  
- sample space, 3
- set subtraction, 3
- simulation, 73
- standard deviation, 30
- standard normal distribution, 76
- stochastic process, 100
  
- uniform distribution, 58
- union, 3
  
- variance, 30
  
- Weibull distribution, 59
  
- zero-sum games, 34

## References

- [KS1991] I. Karatzas and S. Shreve, *Brownian motion and stochastic calculus*, second ed. Springer-Verlag 1991.
- [Laplace1812] P. Laplace, *Théorie analytique des probabilités*, Courcier, Paris, 1812.
- [Malkiel1999] B. Malkiel, *A random walk down Wall Street: the best investment advice for the new century*, W. W. Norton and Company, 1999.
- [Ross1994] S. Ross, *A first course in probability*, 4th Ed. Macmillan College Publishing Co. 1994.
- [Straffin2006] P. Straffin, *Game theory and strategy*, The Mathematical Association of America 1993 (sixth printing 2006).