# Introduction

The analysis of algorithms often requires us to draw upon a body of mathematical tools. Some of these tools are as simple as high-school algebra, but others may be new to you. In Part I, we saw how to manipulate asymptotic notations and solve recurrences. This Appendix is a compendium of several other concepts and methods we use to analyze algorithms. As noted in the introduction to Part I, you may have seen much of the material in this Appendix before having read this book (although the specific notational conventions we use might occasionally differ from those you saw in other books). Hence, you should treat this Appendix as reference material. As in the rest of this book, however, we have included exercises and problems, in order for you to improve your skills in these areas.

Appendix A offers methods for evaluating and bounding summations, which occur frequently in the analysis of algorithms. Many of the formulas in this chapter can be found in any calculus text, but you will find it convenient to have these methods compiled in one place.

Appendix B contains basic definitions and notations for sets, relations, functions, graphs, and trees. This chapter also gives some basic properties of these mathematical objects.

Appendix C begins with elementary principles of counting: permutations, combinations, and the like. The remainder of the chapter contains definitions and properties of basic probability. Most of the algorithms in this book require no probability for their analysis, and thus you can easily omit the latter sections of the chapter on a first reading, even without skimming them. Later, when you encounter a probabilistic analysis that you want to understand better, you will find Appendix C well organized for reference purposes.

# A     Summations

When an algorithm contains an iterative control construct such as a **while** or **for** loop, its running time can be expressed as the sum of the times spent on each execution of the body of the loop. For example, we found in Section 2.2 that the $j$th iteration of insertion sort took time proportional to $j$ in the worst case. By adding up the time spent on each iteration, we obtained the summation (or series)

$$\sum_{j=2}^{n} j .$$

Evaluating this summation yielded a bound of $\Theta(n^2)$ on the worst-case running time of the algorithm. This example indicates the general importance of understanding how to manipulate and bound summations.

Section A.1 lists several basic formulas involving summations. Section A.2 offers useful techniques for bounding summations. The formulas in Section A.1 are given without proof, though proofs for some of them are presented in Section A.2 to illustrate the methods of that section. Most of the other proofs can be found in any calculus text.

## A.1   Summation formulas and properties

Given a sequence $a_1, a_2, \ldots$ of numbers, the finite sum $a_1 + a_2 + \cdots + a_n$, where $n$ is an nonnegative integer, can be written

$$\sum_{k=1}^{n} a_k .$$

If $n = 0$, the value of the summation is defined to be 0. The value of a finite series is always well defined, and its terms can be added in any order.

Given a sequence $a_1, a_2, \ldots$ of numbers, the infinite sum $a_1 + a_2 + \cdots$ can be written

$$\sum_{k=1}^{\infty} a_k \ ,$$

which is interpreted to mean

$$\lim_{n \to \infty} \sum_{k=1}^{n} a_k \ .$$

If the limit does not exist, the series ***diverges***; otherwise, it ***converges***. The terms of a convergent series cannot always be added in any order. We can, however, rearrange the terms of an ***absolutely convergent series***, that is, a series $\sum_{k=1}^{\infty} a_k$ for which the series $\sum_{k=1}^{\infty} |a_k|$ also converges.

## Linearity

For any real number $c$ and any finite sequences $a_1, a_2, \ldots, a_n$ and $b_1, b_2, \ldots, b_n$,

$$\sum_{k=1}^{n} (ca_k + b_k) = c \sum_{k=1}^{n} a_k + \sum_{k=1}^{n} b_k \ .$$

The linearity property is also obeyed by infinite convergent series.

The linearity property can be exploited to manipulate summations incorporating asymptotic notation. For example,

$$\sum_{k=1}^{n} \Theta(f(k)) = \Theta\left(\sum_{k=1}^{n} f(k)\right) \ .$$

In this equation, the $\Theta$-notation on the left-hand side applies to the variable $k$, but on the right-hand side, it applies to $n$. Such manipulations can also be applied to infinite convergent series.

## Arithmetic series

The summation

$$\sum_{k=1}^{n} k = 1 + 2 + \cdots + n \ ,$$

is an ***arithmetic series*** and has the value

$$\sum_{k=1}^{n} k \ = \ \frac{1}{2} n(n+1) \tag{A.1}$$

$$= \ \Theta(n^2) \ . \tag{A.2}$$

### Sums of squares and cubes

We have the following summations of squares and cubes:

$$\sum_{k=0}^{n} k^2 = \frac{n(n+1)(2n+1)}{6} , \tag{A.3}$$

$$\sum_{k=0}^{n} k^3 = \frac{n^2(n+1)^2}{4} . \tag{A.4}$$

### Geometric series

For real $x \neq 1$, the summation

$$\sum_{k=0}^{n} x^k = 1 + x + x^2 + \cdots + x^n$$

is a **geometric** or **exponential series** and has the value

$$\sum_{k=0}^{n} x^k = \frac{x^{n+1} - 1}{x - 1} . \tag{A.5}$$

When the summation is infinite and $|x| < 1$, we have the infinite decreasing geometric series

$$\sum_{k=0}^{\infty} x^k = \frac{1}{1 - x} . \tag{A.6}$$

### Harmonic series

For positive integers $n$, the $n$th **harmonic number** is

$$\begin{aligned}
H_n &= 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \cdots + \frac{1}{n} \\
&= \sum_{k=1}^{n} \frac{1}{k} \\
&= \ln n + O(1) .
\end{aligned} \tag{A.7}$$

(We shall prove this bound in Section A.2.)

### Integrating and differentiating series

Additional formulas can be obtained by integrating or differentiating the formulas above. For example, by differentiating both sides of the infinite geometric series (A.6) and multiplying by $x$, we get

$$\sum_{k=0}^{\infty} kx^k = \frac{x}{(1-x)^2} \tag{A.8}$$

for $|x| < 1$.

**Telescoping series**

For any sequence $a_0, a_1, \ldots, a_n$,

$$\sum_{k=1}^{n} (a_k - a_{k-1}) = a_n - a_0 , \tag{A.9}$$

since each of the terms $a_1, a_2, \ldots, a_{n-1}$ is added in exactly once and subtracted out exactly once. We say that the sum *telescopes*. Similarly,

$$\sum_{k=0}^{n-1} (a_k - a_{k+1}) = a_0 - a_n .$$

As an example of a telescoping sum, consider the series

$$\sum_{k=1}^{n-1} \frac{1}{k(k+1)} .$$

Since we can rewrite each term as

$$\frac{1}{k(k+1)} = \frac{1}{k} - \frac{1}{k+1} ,$$

we get

$$
\begin{aligned}
\sum_{k=1}^{n-1} \frac{1}{k(k+1)} &= \sum_{k=1}^{n-1} \left( \frac{1}{k} - \frac{1}{k+1} \right) \\
&= 1 - \frac{1}{n} .
\end{aligned}
$$

**Products**

The finite product $a_1 a_2 \cdots a_n$ can be written

$$\prod_{k=1}^{n} a_k .$$

If $n = 0$, the value of the product is defined to be 1. We can convert a formula with a product to a formula with a summation by using the identity

$$\lg \left( \prod_{k=1}^{n} a_k \right) = \sum_{k=1}^{n} \lg a_k .$$

**Exercises**

***A.1-1***
Find a simple formula for $\sum_{k=1}^{n}(2k-1)$.

***A.1-2*** ★
Show that $\sum_{k=1}^{n} 1/(2k-1) = \ln(\sqrt{n})+O(1)$ by manipulating the harmonic series.

***A.1-3***
Show that $\sum_{k=0}^{\infty} k^2 x^k = x(1+x)/(1-x)^3$ for $0 < |x| < 1$.

***A.1-4*** ★
Show that $\sum_{k=0}^{\infty}(k-1)/2^k = 0$.

***A.1-5*** ★
Evaluate the sum $\sum_{k=1}^{\infty}(2k+1)x^{2k}$.

***A.1-6***
Prove that $\sum_{k=1}^{n} O(f_k(n)) = O\left(\sum_{k=1}^{n} f_k(n)\right)$ by using the linearity property of summations.

***A.1-7***
Evaluate the product $\prod_{k=1}^{n} 2 \cdot 4^k$.

***A.1-8*** ★
Evaluate the product $\prod_{k=2}^{n}(1 - 1/k^2)$.

---

## A.2   Bounding summations

There are many techniques available for bounding the summations that describe the running times of algorithms. Here are some of the most frequently used methods.

**Mathematical induction**

The most basic way to evaluate a series is to use mathematical induction. As an example, let us prove that the arithmetic series $\sum_{k=1}^{n} k$ evaluates to $\frac{1}{2}n(n+1)$. We can easily verify this for $n = 1$, so we make the inductive assumption that it holds for $n$ and prove that it holds for $n + 1$. We have

$$\sum_{k=1}^{n+1} k = \sum_{k=1}^{n} k + (n+1)$$

$$= \frac{1}{2}n(n+1) + (n+1)$$

$$= \frac{1}{2}(n+1)(n+2) \, .$$

One need not guess the exact value of a summation in order to use mathematical induction. Induction can be used to show a bound as well. As an example, let us prove that the geometric series $\sum_{k=0}^{n} 3^k$ is $O(3^n)$. More specifically, let us prove that $\sum_{k=0}^{n} 3^k \leq c3^n$ for some constant $c$. For the initial condition $n = 0$, we have $\sum_{k=0}^{0} 3^k = 1 \leq c \cdot 1$ as long as $c \geq 1$. Assuming that the bound holds for $n$, let us prove that it holds for $n + 1$. We have

$$\sum_{k=0}^{n+1} 3^k = \sum_{k=0}^{n} 3^k + 3^{n+1}$$

$$\leq c3^n + 3^{n+1}$$

$$= \left( \frac{1}{3} + \frac{1}{c} \right) c3^{n+1}$$

$$\leq c3^{n+1}$$

as long as $(1/3 + 1/c) \leq 1$ or, equivalently, $c \geq 3/2$. Thus, $\sum_{k=0}^{n} 3^k = O(3^n)$, as we wished to show.

We have to be careful when we use asymptotic notation to prove bounds by induction. Consider the following fallacious proof that $\sum_{k=1}^{n} k = O(n)$. Certainly, $\sum_{k=1}^{1} k = O(1)$. Assuming the bound for $n$, we now prove it for $n + 1$:

$$\sum_{k=1}^{n+1} k = \sum_{k=1}^{n} k + (n+1)$$

$$= O(n) + (n+1) \qquad \Longleftarrow wrong!!$$

$$= O(n+1) \, .$$

The bug in the argument is that the "constant" hidden by the "big-oh" grows with $n$ and thus is not constant. We have not shown that the same constant works for *all n*.

### Bounding the terms

Sometimes, a good upper bound on a series can be obtained by bounding each term of the series, and it often suffices to use the largest term to bound the others. For example, a quick upper bound on the arithmetic series (A.1) is

$$\sum_{k=1}^{n} k \leq \sum_{k=1}^{n} n$$

$$= n^2 \, .$$

In general, for a series $\sum_{k=1}^{n} a_k$, if we let $a_{\max} = \max_{1 \le k \le n} a_k$, then

$$\sum_{k=1}^{n} a_k \le n a_{\max} .$$

The technique of bounding each term in a series by the largest term is a weak method when the series can in fact be bounded by a geometric series. Given the series $\sum_{k=0}^{n} a_k$, suppose that $a_{k+1}/a_k \le r$ for all $k \ge 0$, where $0 < r < 1$ is a constant. The sum can be bounded by an infinite decreasing geometric series, since $a_k \le a_0 r^k$, and thus

$$
\begin{aligned}
\sum_{k=0}^{n} a_k &\le \sum_{k=0}^{\infty} a_0 r^k \\
&= a_0 \sum_{k=0}^{\infty} r^k \\
&= a_0 \frac{1}{1-r} .
\end{aligned}
$$

We can apply this method to bound the summation $\sum_{k=1}^{\infty} (k/3^k)$. In order to start the summation at $k = 0$, we rewrite it as $\sum_{k=0}^{\infty} ((k+1)/3^{k+1})$. The first term $(a_0)$ is $1/3$, and the ratio $(r)$ of consecutive terms is

$$
\begin{aligned}
\frac{(k+2)/3^{k+2}}{(k+1)/3^{k+1}} &= \frac{1}{3} \cdot \frac{k+2}{k+1} \\
&\le \frac{2}{3}
\end{aligned}
$$

for all $k \ge 0$. Thus, we have

$$
\begin{aligned}
\sum_{k=1}^{\infty} \frac{k}{3^k} &= \sum_{k=0}^{\infty} \frac{k+1}{3^{k+1}} \\
&\le \frac{1}{3} \cdot \frac{1}{1 - 2/3} \\
&= 1 .
\end{aligned}
$$

A common bug in applying this method is to show that the ratio of consecutive terms is less than 1 and then to assume that the summation is bounded by a geometric series. An example is the infinite harmonic series, which diverges since

$$
\begin{aligned}
\sum_{k=1}^{\infty} \frac{1}{k} &= \lim_{n \to \infty} \sum_{k=1}^{n} \frac{1}{k} \\
&= \lim_{n \to \infty} \Theta(\lg n) \\
&= \infty .
\end{aligned}
$$

The ratio of the $(k+1)$st and $k$th terms in this series is $k/(k+1) < 1$, but the series is not bounded by a decreasing geometric series. To bound a series by a geometric series, one must show that there is an $r < 1$, which is a *constant*, such that the ratio of all pairs of consecutive terms never exceeds $r$. In the harmonic series, no such $r$ exists because the ratio becomes arbitrarily close to 1.

## Splitting summations

One way to obtain bounds on a difficult summation is to express the series as the sum of two or more series by partitioning the range of the index and then to bound each of the resulting series. For example, suppose we try to find a lower bound on the arithmetic series $\sum_{k=1}^{n} k$, which has already been shown to have an upper bound of $n^2$. We might attempt to bound each term in the summation by the smallest term, but since that term is 1, we get a lower bound of $n$ for the summation—far off from our upper bound of $n^2$.

We can obtain a better lower bound by first splitting the summation. Assume for convenience that $n$ is even. We have

$$
\begin{aligned}
\sum_{k=1}^{n} k &= \sum_{k=1}^{n/2} k + \sum_{k=n/2+1}^{n} k \\
&\geq \sum_{k=1}^{n/2} 0 + \sum_{k=n/2+1}^{n} (n/2) \\
&= (n/2)^2 \\
&= \Omega(n^2) ,
\end{aligned}
$$

which is an asymptotically tight bound, since $\sum_{k=1}^{n} k = O(n^2)$.

For a summation arising from the analysis of an algorithm, we can often split the summation and ignore a constant number of the initial terms. Generally, this technique applies when each term $a_k$ in a summation $\sum_{k=0}^{n} a_k$ is independent of $n$. Then for any constant $k_0 > 0$, we can write

$$
\begin{aligned}
\sum_{k=0}^{n} a_k &= \sum_{k=0}^{k_0-1} a_k + \sum_{k=k_0}^{n} a_k \\
&= \Theta(1) + \sum_{k=k_0}^{n} a_k ,
\end{aligned}
$$

since the initial terms of the summation are all constant and there are a constant number of them. We can then use other methods to bound $\sum_{k=k_0}^{n} a_k$. This technique applies to infinite summations as well. For example, to find an asymptotic upper bound on

$$\sum_{k=0}^{\infty} \frac{k^2}{2^k} \ ,$$

we observe that the ratio of consecutive terms is

$$\frac{(k+1)^2/2^{k+1}}{k^2/2^k} \ = \ \frac{(k+1)^2}{2k^2}$$

$$\leq \ \frac{8}{9}$$

if $k \geq 3$. Thus, the summation can be split into

$$\sum_{k=0}^{\infty} \frac{k^2}{2^k} \ = \ \sum_{k=0}^{2} \frac{k^2}{2^k} + \sum_{k=3}^{\infty} \frac{k^2}{2^k}$$

$$\leq \ \sum_{k=0}^{2} \frac{k^2}{2^k} + \frac{9}{8} \sum_{k=0}^{\infty} \left(\frac{8}{9}\right)^k$$

$$= \ O(1) \ ,$$

since the first summation has a constant number of terms and the second summation is a decreasing geometric series.

The technique of splitting summations can be used to determine asymptotic bounds in much more difficult situations. For example, we can obtain a bound of $O(\lg n)$ on the harmonic series (A.7):

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

The idea is to split the range 1 to $n$ into $\lfloor \lg n \rfloor$ pieces and upper-bound the contribution of each piece by 1. Each piece consists of the terms starting at $1/2^i$ and going up to but not including $1/2^{i+1}$, giving

$$\sum_{k=1}^{n} \frac{1}{k} \ \leq \ \sum_{i=0}^{\lfloor \lg n \rfloor} \sum_{j=0}^{2^i-1} \frac{1}{2^i + j}$$

$$\leq \ \sum_{i=0}^{\lfloor \lg n \rfloor} \sum_{j=0}^{2^i-1} \frac{1}{2^i}$$

$$\leq \ \sum_{i=0}^{\lfloor \lg n \rfloor} 1$$

$$\leq \ \lg n + 1 \ . \tag{A.10}$$

### Approximation by integrals

When a summation can be expressed as $\sum_{k=m}^{n} f(k)$, where $f(k)$ is a monotonically increasing function, we can approximate it by integrals:

$$\int_{m-1}^{n} f(x)\,dx \le \sum_{k=m}^{n} f(k) \le \int_{m}^{n+1} f(x)\,dx \ . \tag{A.11}$$

The justification for this approximation is shown in Figure A.1. The summation is represented as the area of the rectangles in the figure, and the integral is the shaded region under the curve. When $f(k)$ is a monotonically decreasing function, we can use a similar method to provide the bounds

$$\int_{m}^{n+1} f(x)\,dx \le \sum_{k=m}^{n} f(k) \le \int_{m-1}^{n} f(x)\,dx \ . \tag{A.12}$$

The integral approximation (A.12) gives a tight estimate for the $n$th harmonic number. For a lower bound, we obtain

$$\begin{aligned}
\sum_{k=1}^{n} \frac{1}{k} &\ge \int_{1}^{n+1} \frac{dx}{x} \\
&= \ln(n+1) \ . 
\end{aligned} \tag{A.13}$$

For the upper bound, we derive the inequality

$$\begin{aligned}
\sum_{k=2}^{n} \frac{1}{k} &\le \int_{1}^{n} \frac{dx}{x} \\
&= \ln n \ ,
\end{aligned}$$

which yields the bound

$$\sum_{k=1}^{n} \frac{1}{k} \le \ln n + 1 \ . \tag{A.14}$$

### Exercises

***A.2-1***
Show that $\sum_{k=1}^{n} 1/k^2$ is bounded above by a constant.

***A.2-2***
Find an asymptotic upper bound on the summation

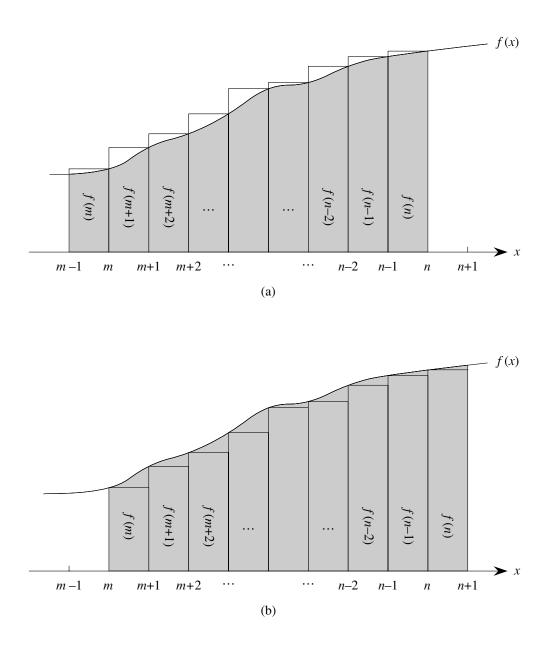$$\sum_{k=0}^{\lfloor \lg n \rfloor} \lceil n/2^k \rceil \ .$$

**Figure A.1**   Approximation of $\sum_{k=m}^{n} f(k)$ by integrals. The area of each rectangle is shown within the rectangle, and the total rectangle area represents the value of the summation. The integral is represented by the shaded area under the curve. By comparing areas in **(a)**, we get $\int_{m-1}^{n} f(x)\,dx \ \leq \ \sum_{k=m}^{n} f(k)$, and then by shifting the rectangles one unit to the right, we get $\sum_{k=m}^{n} f(k) \leq \int_{m}^{n+1} f(x)\,dx$ in **(b)**.

### A.2-3
Show that the $n$th harmonic number is $\Omega(\lg n)$ by splitting the summation.

### A.2-4
Approximate $\sum_{k=1}^{n} k^3$ with an integral.

### A.2-5
Why didn't we use the integral approximation (A.12) directly on $\sum_{k=1}^{n} 1/k$ to obtain an upper bound on the $n$th harmonic number?

## Problems

### A-1 *Bounding summations*
Give asymptotically tight bounds on the following summations. Assume that $r \geq 0$ and $s \geq 0$ are constants.

*a.* $\displaystyle\sum_{k=1}^{n} k^r$.

*b.* $\displaystyle\sum_{k=1}^{n} \lg^s k$.

*c.* $\displaystyle\sum_{k=1}^{n} k^r \lg^s k$.

## Chapter notes

Knuth [182] is an excellent reference for the material presented in this chapter. Basic properties of series can be found in any good calculus book, such as Apostol [18] or Thomas and Finney [296].

# C          Counting and Probability

This chapter reviews elementary combinatorics and probability theory. If you have a good background in these areas, you may want to skim the beginning of the chapter lightly and concentrate on the later sections. Most of the chapters do not require probability, but for some chapters it is essential.

Section C.1 reviews elementary results in counting theory, including standard formulas for counting permutations and combinations. The axioms of probability and basic facts concerning probability distributions are presented in Section C.2. Random variables are introduced in Section C.3, along with the properties of expectation and variance. Section C.4 investigates the geometric and binomial distributions that arise from studying Bernoulli trials. The study of the binomial distribution is continued in Section C.5, an advanced discussion of the "tails" of the distribution.

## C.1   Counting

Counting theory tries to answer the question "How many?" without actually enumerating how many. For example, we might ask, "How many different $n$-bit numbers are there?" or "How many orderings of $n$ distinct elements are there?" In this section, we review the elements of counting theory. Since some of the material assumes a basic understanding of sets, the reader is advised to start by reviewing the material in Section B.1.

### Rules of sum and product

A set of items that we wish to count can sometimes be expressed as a union of disjoint sets or as a Cartesian product of sets.

The *rule of sum* says that the number of ways to choose an element from one of two *disjoint* sets is the sum of the cardinalities of the sets. That is, if $A$ and $B$ are two finite sets with no members in common, then $|A \cup B| = |A| + |B|$, which

follows from equation (B.3). For example, each position on a car's license plate is a letter or a digit. The number of possibilities for each position is therefore $26 + 10 = 36$, since there are 26 choices if it is a letter and 10 choices if it is a digit.

The **rule of product** says that the number of ways to choose an ordered pair is the number of ways to choose the first element times the number of ways to choose the second element. That is, if $A$ and $B$ are two finite sets, then $|A \times B| = |A| \cdot |B|$, which is simply equation (B.4). For example, if an ice-cream parlor offers 28 flavors of ice cream and 4 toppings, the number of possible sundaes with one scoop of ice cream and one topping is $28 \cdot 4 = 112$.

## Strings

A **string** over a finite set $S$ is a sequence of elements of $S$. For example, there are 8 binary strings of length 3:

$$000, 001, 010, 011, 100, 101, 110, 111 .$$

We sometimes call a string of length $k$ a **$k$-string**. A **substring** $s'$ of a string $s$ is an ordered sequence of consecutive elements of $s$. A **$k$-substring** of a string is a substring of length $k$. For example, 010 is a 3-substring of 01101001 (the 3-substring that begins in position 4), but 111 is not a substring of 01101001.

A $k$-string over a set $S$ can be viewed as an element of the Cartesian product $S^k$ of $k$-tuples; thus, there are $|S|^k$ strings of length $k$. For example, the number of binary $k$-strings is $2^k$. Intuitively, to construct a $k$-string over an $n$-set, we have $n$ ways to pick the first element; for each of these choices, we have $n$ ways to pick the second element; and so forth $k$ times. This construction leads to the $k$-fold product $n \cdot n \cdots n = n^k$ as the number of $k$-strings.

## Permutations

A **permutation** of a finite set $S$ is an ordered sequence of all the elements of $S$, with each element appearing exactly once. For example, if $S = \{a, b, c\}$, there are 6 permutations of $S$:

$$abc, acb, bac, bca, cab, cba .$$

There are $n!$ permutations of a set of $n$ elements, since the first element of the sequence can be chosen in $n$ ways, the second in $n - 1$ ways, the third in $n - 2$ ways, and so on.

A **$k$-permutation** of $S$ is an ordered sequence of $k$ elements of $S$, with no element appearing more than once in the sequence. (Thus, an ordinary permutation is just an $n$-permutation of an $n$-set.) The twelve 2-permutations of the set $\{a, b, c, d\}$ are

$$ab, ac, ad, ba, bc, bd, ca, cb, cd, da, db, dc .$$

The number of $k$-permutations of an $n$-set is

$$n(n - 1)(n - 2) \cdots (n - k + 1) = \frac{n!}{(n - k)!} \ , \tag{C.1}$$

since there are $n$ ways of choosing the first element, $n - 1$ ways of choosing the second element, and so on until $k$ elements are selected, the last being a selection from $n - k + 1$ elements.

## Combinations

A **$k$-combination** of an $n$-set $S$ is simply a $k$-subset of $S$. For example, there are six 2-combinations of the 4-set $\{a, b, c, d\}$:

$$ab, ac, ad, bc, bd, cd \ .$$

(Here we use the shorthand of denoting the 2-set $\{a, b\}$ by $ab$, and so on.) We can construct a $k$-combination of an $n$-set by choosing $k$ distinct (different) elements from the $n$-set.

The number of $k$-combinations of an $n$-set can be expressed in terms of the number of $k$-permutations of an $n$-set. For every $k$-combination, there are exactly $k!$ permutations of its elements, each of which is a distinct $k$-permutation of the $n$-set. Thus, the number of $k$-combinations of an $n$-set is the number of $k$-permutations divided by $k!$; from equation (C.1), this quantity is

$$\frac{n!}{k! \, (n - k)!} \ . \tag{C.2}$$

For $k = 0$, this formula tells us that the number of ways to choose 0 elements from an $n$-set is 1 (not 0), since $0! = 1$.

## Binomial coefficients

We use the notation $\binom{n}{k}$ (read "$n$ choose $k$") to denote the number of $k$-combinations of an $n$-set. From equation (C.2), we have

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

This formula is symmetric in $k$ and $n - k$:

$$\binom{n}{k} = \binom{n}{n - k} \ . \tag{C.3}$$

These numbers are also known as **binomial coefficients**, due to their appearance in the **binomial expansion**:

$$(x + y)^n = \sum_{k=0}^{n} \binom{n}{k} x^k y^{n-k} . \tag{C.4}$$

A special case of the binomial expansion occurs when $x = y = 1$:

$$2^n = \sum_{k=0}^{n} \binom{n}{k} .$$

This formula corresponds to counting the $2^n$ binary $n$-strings by the number of 1's they contain: there are $\binom{n}{k}$ binary $n$-strings containing exactly $k$ 1's, since there are $\binom{n}{k}$ ways to choose $k$ out of the $n$ positions in which to place the 1's.

There are many identities involving binomial coefficients. The exercises at the end of this section give you the opportunity to prove a few.

**Binomial bounds**

We sometimes need to bound the size of a binomial coefficient. For $1 \le k \le n$, we have the lower bound

$$\begin{aligned}
\binom{n}{k} &= \frac{n(n-1)\cdots(n-k+1)}{k(k-1)\cdots 1} \\
&= \left(\frac{n}{k}\right)\left(\frac{n-1}{k-1}\right)\cdots\left(\frac{n-k+1}{1}\right) \\
&\ge \left(\frac{n}{k}\right)^k .
\end{aligned}$$

Taking advantage of the inequality $k! \ge (k/e)^k$ derived from Stirling's approximation (3.17), we obtain the upper bounds

$$\begin{aligned}
\binom{n}{k} &= \frac{n(n-1)\cdots(n-k+1)}{k(k-1)\cdots 1} \\
&\le \frac{n^k}{k!} \\
&\le \left(\frac{en}{k}\right)^k . \tag{C.5}
\end{aligned}$$

For all $0 \le k \le n$, we can use induction (see Exercise C.1-12) to prove the bound

$$\binom{n}{k} \le \frac{n^n}{k^k(n-k)^{n-k}} , \tag{C.6}$$

where for convenience we assume that $0^0 = 1$. For $k = \lambda n$, where $0 \le \lambda \le 1$, this bound can be rewritten as

$$\binom{n}{\lambda n} \le \frac{n^n}{(\lambda n)^{\lambda n}((1-\lambda)n)^{(1-\lambda)n}}$$

$$= \left( \left( \frac{1}{\lambda} \right)^{\lambda} \left( \frac{1}{1 - \lambda} \right)^{1-\lambda} \right)^{n}$$
$$= 2^{n H(\lambda)} ,$$

where

$$H(\lambda) = -\lambda \lg \lambda - (1 - \lambda) \lg(1 - \lambda) \tag{C.7}$$

is the *(binary) entropy function* and where, for convenience, we assume that $0 \lg 0 = 0$, so that $H(0) = H(1) = 0$.

### Exercises

#### *C.1-1*
How many $k$-substrings does an $n$-string have? (Consider identical $k$-substrings at different positions as different.) How many substrings does an $n$-string have in total?

#### *C.1-2*
An $n$-input, $m$-output *boolean function* is a function from $\{\text{TRUE, FALSE}\}^{n}$ to $\{\text{TRUE, FALSE}\}^{m}$. How many $n$-input, 1-output boolean functions are there? How many $n$-input, $m$-output boolean functions are there?

#### *C.1-3*
In how many ways can $n$ professors sit around a circular conference table? Consider two seatings to be the same if one can be rotated to form the other.

#### *C.1-4*
How many ways are there to choose from the set $\{1, 2, \ldots, 100\}$ three distinct numbers so that their sum is even?

#### *C.1-5*
Prove the identity

$$\binom{n}{k} = \frac{n}{k} \binom{n-1}{k-1} \tag{C.8}$$

for $0 < k \leq n$.

#### *C.1-6*
Prove the identity

$$\binom{n}{k} = \frac{n}{n-k} \binom{n-1}{k}$$

for $0 \leq k < n$.

### C.1-7

To choose $k$ objects from $n$, you can make one of the objects distinguished and consider whether the distinguished object is chosen. Use this approach to prove that

$$\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}.$$

### C.1-8

Using the result of Exercise C.1-7, make a table for $n = 0, 1, \ldots, 6$ and $0 \le k \le n$ of the binomial coefficients $\binom{n}{k}$ with $\binom{0}{0}$ at the top, $\binom{1}{0}$ and $\binom{1}{1}$ on the next line, and so forth. Such a table of binomial coefficients is called *Pascal's triangle*.

### C.1-9

Prove that

$$\sum_{i=1}^{n} i = \binom{n+1}{2}.$$

### C.1-10

Show that for any $n \ge 0$ and $0 \le k \le n$, the maximum value of $\binom{n}{k}$ is achieved when $k = \lfloor n/2 \rfloor$ or $k = \lceil n/2 \rceil$.

### C.1-11 ★

Argue that for any $n \ge 0$, $j \ge 0$, $k \ge 0$, and $j + k \le n$,

$$\binom{n}{j+k} \le \binom{n}{j}\binom{n-j}{k}. \tag{C.9}$$

Provide both an algebraic proof and an argument based on a method for choosing $j + k$ items out of $n$. Give an example in which equality does not hold.

### C.1-12 ★

Use induction on $k \le n/2$ to prove inequality (C.6), and use equation (C.3) to extend it to all $k \le n$.

### C.1-13 ★

Use Stirling's approximation to prove that

$$\binom{2n}{n} = \frac{2^{2n}}{\sqrt{\pi n}}(1 + O(1/n)). \tag{C.10}$$

### C.1-14 ★

By differentiating the entropy function $H(\lambda)$, show that it achieves its maximum value at $\lambda = 1/2$. What is $H(1/2)$?

### C.1-15 ⋆

Show that for any integer $n \geq 0$,

$$\sum_{k=0}^{n} \binom{n}{k} k = n2^{n-1} . \tag{C.11}$$

## C.2    Probability

Probability is an essential tool for the design and analysis of probabilistic and randomized algorithms. This section reviews basic probability theory.

We define probability in terms of a ***sample space*** $S$, which is a set whose elements are called ***elementary events***. Each elementary event can be viewed as a possible outcome of an experiment. For the experiment of flipping two distinguishable coins, we can view the sample space as consisting of the set of all possible 2-strings over {H, T}:

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

An ***event*** is a subset[1] of the sample space $S$. For example, in the experiment of flipping two coins, the event of obtaining one head and one tail is {HT, TH}. The event $S$ is called the ***certain event***, and the event $\emptyset$ is called the ***null event***. We say that two events $A$ and $B$ are ***mutually exclusive*** if $A \cap B = \emptyset$. We sometimes treat an elementary event $s \in S$ as the event $\{s\}$. By definition, all elementary events are mutually exclusive.

### Axioms of probability

A ***probability distribution*** $\Pr\{\}$ on a sample space $S$ is a mapping from events of $S$ to real numbers such that the following ***probability axioms*** are satisfied:

1.  $\Pr\{A\} \geq 0$ for any event $A$.
2.  $\Pr\{S\} = 1$.

---

[1] For a general probability distribution, there may be some subsets of the sample space $S$ that are not considered to be events. This situation usually arises when the sample space is uncountably infinite. The main requirement is that the set of events of a sample space be closed under the operations of taking the complement of an event, forming the union of a finite or countable number of events, and taking the intersection of a finite or countable number of events. Most of the probability distributions we shall see are over finite or countable sample spaces, and we shall generally consider all subsets of a sample space to be events. A notable exception is the continuous uniform probability distribution, which will be presented shortly.

3. $\Pr\{A \cup B\} = \Pr\{A\} + \Pr\{B\}$ for any two mutually exclusive events $A$ and $B$. More generally, for any (finite or countably infinite) sequence of events $A_1, A_2, \ldots$ that are pairwise mutually exclusive,

$$\Pr\left\{\bigcup_i A_i\right\} = \sum_i \Pr\{A_i\} \ .$$

We call $\Pr\{A\}$ the **probability** of the event $A$. We note here that axiom 2 is a normalization requirement: there is really nothing fundamental about choosing 1 as the probability of the certain event, except that it is natural and convenient.

Several results follow immediately from these axioms and basic set theory (see Section B.1). The null event $\emptyset$ has probability $\Pr\{\emptyset\} = 0$. If $A \subseteq B$, then $\Pr\{A\} \leq \Pr\{B\}$. Using $\overline{A}$ to denote the event $S - A$ (the **complement** of $A$), we have $\Pr\{\overline{A}\} = 1 - \Pr\{A\}$. For any two events $A$ and $B$,

$$\Pr\{A \cup B\} = \Pr\{A\} + \Pr\{B\} - \Pr\{A \cap B\} \tag{C.12}$$
$$\leq \Pr\{A\} + \Pr\{B\} \ . \tag{C.13}$$

In our coin-flipping example, suppose that each of the four elementary events has probability $1/4$. Then the probability of getting at least one head is

$$\Pr\{\text{HH, HT, TH}\} = \Pr\{\text{HH}\} + \Pr\{\text{HT}\} + \Pr\{\text{TH}\}$$
$$= 3/4 \ .$$

Alternatively, since the probability of getting strictly less than one head is $\Pr\{\text{TT}\} = 1/4$, the probability of getting at least one head is $1 - 1/4 = 3/4$.

### Discrete probability distributions

A probability distribution is **discrete** if it is defined over a finite or countably infinite sample space. Let $S$ be the sample space. Then for any event $A$,

$$\Pr\{A\} = \sum_{s \in A} \Pr\{s\} \ ,$$

since elementary events, specifically those in $A$, are mutually exclusive. If $S$ is finite and every elementary event $s \in S$ has probability

$$\Pr\{s\} = 1/|S| \ ,$$

then we have the **uniform probability distribution** on $S$. In such a case the experiment is often described as "picking an element of $S$ at random."

As an example, consider the process of flipping a **fair coin**, one for which the probability of obtaining a head is the same as the probability of obtaining a tail, that is, $1/2$. If we flip the coin $n$ times, we have the uniform probability distribution

defined on the sample space $S = \{H, T\}^n$, a set of size $2^n$. Each elementary event in $S$ can be represented as a string of length $n$ over $\{H, T\}$, and each occurs with probability $1/2^n$. The event

$A = \{$exactly $k$ heads and exactly $n - k$ tails occur$\}$

is a subset of $S$ of size $|A| = \binom{n}{k}$, since there are $\binom{n}{k}$ strings of length $n$ over $\{H, T\}$ that contain exactly $k$ H's. The probability of event $A$ is thus $\Pr\{A\} = \binom{n}{k}/2^n$.

### Continuous uniform probability distribution

The continuous uniform probability distribution is an example of a probability distribution in which not all subsets of the sample space are considered to be events. The continuous uniform probability distribution is defined over a closed interval $[a, b]$ of the reals, where $a < b$. Intuitively, we want each point in the interval $[a, b]$ to be "equally likely." There is an uncountable number of points, however, so if we give all points the same finite, positive probability, we cannot simultaneously satisfy axioms 2 and 3. For this reason, we would like to associate a probability only with *some* of the subsets of $S$ in such a way that the axioms are satisfied for these events.

For any closed interval $[c, d]$, where $a \leq c \leq d \leq b$, the **continuous uniform probability distribution** defines the probability of the event $[c, d]$ to be

$$\Pr\{[c, d]\} = \frac{d - c}{b - a} \ .$$

Note that for any point $x = [x, x]$, the probability of $x$ is 0. If we remove the endpoints of an interval $[c, d]$, we obtain the open interval $(c, d)$. Since $[c, d] = [c, c] \cup (c, d) \cup [d, d]$, axiom 3 gives us $\Pr\{[c, d]\} = \Pr\{(c, d)\}$. Generally, the set of events for the continuous uniform probability distribution is any subset of the sample space $[a, b]$ that can be obtained by a finite or countable union of open and closed intervals.

### Conditional probability and independence

Sometimes we have some prior partial knowledge about the outcome of an experiment. For example, suppose that a friend has flipped two fair coins and has told you that at least one of the coins showed a head. What is the probability that both coins are heads? The information given eliminates the possibility of two tails. The three remaining elementary events are equally likely, so we infer that each occurs with probability 1/3. Since only one of these elementary events shows two heads, the answer to our question is 1/3.

Conditional probability formalizes the notion of having prior partial knowledge of the outcome of an experiment. The ***conditional probability*** of an event $A$ given that another event $B$ occurs is defined to be

$$\Pr\{A \mid B\} = \frac{\Pr\{A \cap B\}}{\Pr\{B\}} \tag{C.14}$$

whenever $\Pr\{B\} \neq 0$. (We read "$\Pr\{A \mid B\}$" as "the probability of $A$ given $B$.") Intuitively, since we are given that event $B$ occurs, the event that $A$ also occurs is $A \cap B$. That is, $A \cap B$ is the set of outcomes in which both $A$ and $B$ occur. Since the outcome is one of the elementary events in $B$, we normalize the probabilities of all the elementary events in $B$ by dividing them by $\Pr\{B\}$, so that they sum to 1. The conditional probability of $A$ given $B$ is, therefore, the ratio of the probability of event $A \cap B$ to the probability of event $B$. In the example above, $A$ is the event that both coins are heads, and $B$ is the event that at least one coin is a head. Thus, $\Pr\{A \mid B\} = (1/4)/(3/4) = 1/3$.

Two events are ***independent*** if

$$\Pr\{A \cap B\} = \Pr\{A\}\Pr\{B\} \ , \tag{C.15}$$

which is equivalent, if $\Pr\{B\} \neq 0$, to the condition

$$\Pr\{A \mid B\} = \Pr\{A\} \ .$$

For example, suppose that two fair coins are flipped and that the outcomes are independent. Then the probability of two heads is $(1/2)(1/2) = 1/4$. Now suppose that one event is that the first coin comes up heads and the other event is that the coins come up differently. Each of these events occurs with probability $1/2$, and the probability that both events occur is $1/4$; thus, according to the definition of independence, the events are independent—even though one might think that both events depend on the first coin. Finally, suppose that the coins are welded together so that they both fall heads or both fall tails and that the two possibilities are equally likely. Then the probability that each coin comes up heads is $1/2$, but the probability that they both come up heads is $1/2 \neq (1/2)(1/2)$. Consequently, the event that one comes up heads and the event that the other comes up heads are not independent.

A collection $A_1, A_2, \ldots, A_n$ of events is said to be ***pairwise independent*** if

$$\Pr\{A_i \cap A_j\} = \Pr\{A_i\}\Pr\{A_j\}$$

for all $1 \leq i < j \leq n$. We say that the events of the collection are ***(mutually) independent*** if every $k$-subset $A_{i_1}, A_{i_2}, \ldots, A_{i_k}$ of the collection, where $2 \leq k \leq n$ and $1 \leq i_1 < i_2 < \cdots < i_k \leq n$, satisfies

$$\Pr\{A_{i_1} \cap A_{i_2} \cap \cdots \cap A_{i_k}\} = \Pr\{A_{i_1}\}\Pr\{A_{i_2}\} \cdots \Pr\{A_{i_k}\} \ .$$

For example, suppose we flip two fair coins. Let $A_1$ be the event that the first coin is heads, let $A_2$ be the event that the second coin is heads, and let $A_3$ be the event that the two coins are different. We have

$$
\begin{aligned}
\Pr\{A_1\} &= 1/2, \\
\Pr\{A_2\} &= 1/2, \\
\Pr\{A_3\} &= 1/2, \\
\Pr\{A_1 \cap A_2\} &= 1/4, \\
\Pr\{A_1 \cap A_3\} &= 1/4, \\
\Pr\{A_2 \cap A_3\} &= 1/4, \\
\Pr\{A_1 \cap A_2 \cap A_3\} &= 0.
\end{aligned}
$$

Since for $1 \leq i < j \leq 3$, we have $\Pr\{A_i \cap A_j\} = \Pr\{A_i\}\Pr\{A_j\} = 1/4$, the events $A_1$, $A_2$, and $A_3$ are pairwise independent. The events are not mutually independent, however, because $\Pr\{A_1 \cap A_2 \cap A_3\} = 0$ and $\Pr\{A_1\}\Pr\{A_2\}\Pr\{A_3\} = 1/8 \neq 0$.

**Bayes's theorem**

From the definition of conditional probability (C.14) and the commutative law $A \cap B = B \cap A$, it follows that for two events $A$ and $B$, each with nonzero probability,

$$
\begin{aligned}
\Pr\{A \cap B\} &= \Pr\{B\}\Pr\{A \mid B\} \\
&= \Pr\{A\}\Pr\{B \mid A\}.
\end{aligned}
\tag{C.16}
$$

Solving for $\Pr\{A \mid B\}$, we obtain

$$
\Pr\{A \mid B\} = \frac{\Pr\{A\}\Pr\{B \mid A\}}{\Pr\{B\}},
\tag{C.17}
$$

which is known as ***Bayes's theorem***. The denominator $\Pr\{B\}$ is a normalizing constant that we can reexpress as follows. Since $B = (B \cap A) \cup (B \cap \overline{A})$ and $B \cap A$ and $B \cap \overline{A}$ are mutually exclusive events,

$$
\begin{aligned}
\Pr\{B\} &= \Pr\{B \cap A\} + \Pr\{B \cap \overline{A}\} \\
&= \Pr\{A\}\Pr\{B \mid A\} + \Pr\{\overline{A}\}\Pr\{B \mid \overline{A}\}.
\end{aligned}
$$

Substituting into equation (C.17), we obtain an equivalent form of Bayes's theorem:

$$
\Pr\{A \mid B\} = \frac{\Pr\{A\}\Pr\{B \mid A\}}{\Pr\{A\}\Pr\{B \mid A\} + \Pr\{\overline{A}\}\Pr\{B \mid \overline{A}\}}.
$$

Bayes's theorem can simplify the computing of conditional probabilities. For example, suppose that we have a fair coin and a biased coin that always comes up heads. We run an experiment consisting of three independent events: one of the two coins is chosen at random, the coin is flipped once, and then it is flipped again. Suppose that the chosen coin comes up heads both times. What is the probability that it is biased?

We solve this problem using Bayes's theorem. Let $A$ be the event that the biased coin is chosen, and let $B$ be the event that the coin comes up heads both times. We wish to determine $\Pr\{A \mid B\}$. We have $\Pr\{A\} = 1/2$, $\Pr\{B \mid A\} = 1$, $\Pr\{\overline{A}\} = 1/2$, and $\Pr\{B \mid \overline{A}\} = 1/4$; hence,

$$
\begin{aligned}
\Pr\{A \mid B\} &= \frac{(1/2) \cdot 1}{(1/2) \cdot 1 + (1/2) \cdot (1/4)} \\
&= 4/5 \ .
\end{aligned}
$$

**Exercises**

***C.2-1***
Prove ***Boole's inequality***: For any finite or countably infinite sequence of events $A_1, A_2, \ldots,$

$$
\Pr\{A_1 \cup A_2 \cup \cdots\} \leq \Pr\{A_1\} + \Pr\{A_2\} + \cdots \ . \tag{C.18}
$$

***C.2-2***
Professor Rosencrantz flips a fair coin once. Professor Guildenstern flips a fair coin twice. What is the probability that Professor Rosencrantz obtains more heads than Professor Guildenstern?

***C.2-3***
A deck of 10 cards, each bearing a distinct number from 1 to 10, is shuffled to mix the cards thoroughly. Three cards are removed one at a time from the deck. What is the probability that the three cards are selected in sorted (increasing) order?

***C.2-4*** ★
Describe a procedure that takes as input two integers $a$ and $b$ such that $0 < a < b$ and, using fair coin flips, produces as output heads with probability $a/b$ and tails with probability $(b - a)/b$. Give a bound on the expected number of coin flips, which should be $O(1)$. (*Hint:* Represent $a/b$ in binary.)

***C.2-5***
Prove that

$$
\Pr\{A \mid B\} + \Pr\{\overline{A} \mid B\} = 1 \ .
$$

### C.2-6

Prove that for any collection of events $A_1, A_2, \ldots, A_n$,

$$\Pr\{A_1 \cap A_2 \cap \cdots \cap A_n\} = \Pr\{A_1\} \cdot \Pr\{A_2 \mid A_1\} \cdot \Pr\{A_3 \mid A_1 \cap A_2\} \cdots$$
$$\Pr\{A_n \mid A_1 \cap A_2 \cap \cdots \cap A_{n-1}\} \ .$$

### C.2-7 ★

Show how to construct a set of $n$ events that are pairwise independent but such that no subset of $k > 2$ of them is mutually independent.

### C.2-8 ★

Two events $A$ and $B$ are ***conditionally independent***, given $C$, if

$$\Pr\{A \cap B \mid C\} = \Pr\{A \mid C\} \cdot \Pr\{B \mid C\} \ .$$

Give a simple but nontrivial example of two events that are not independent but are conditionally independent given a third event.

### C.2-9 ★

You are a contestant in a game show in which a prize is hidden behind one of three curtains. You will win the prize if you select the correct curtain. After you have picked one curtain but before the curtain is lifted, the emcee lifts one of the other curtains, knowing that it will reveal an empty stage, and asks if you would like to switch from your current selection to the remaining curtain. How would your chances change if you switch?

### C.2-10 ★

A prison warden has randomly picked one prisoner among three to go free. The other two will be executed. The guard knows which one will go free but is forbidden to give any prisoner information regarding his status. Let us call the prisoners $X, Y$, and $Z$. Prisoner $X$ asks the guard privately which of $Y$ or $Z$ will be executed, arguing that since he already knows that at least one of them must die, the guard won't be revealing any information about his own status. The guard tells $X$ that $Y$ is to be executed. Prisoner $X$ feels happier now, since he figures that either he or prisoner $Z$ will go free, which means that his probability of going free is now 1/2. Is he right, or are his chances still 1/3? Explain.

## C.3   Discrete random variables

A ***(discrete) random variable*** $X$ is a function from a finite or countably infinite sample space $S$ to the real numbers. It associates a real number with each possible

outcome of an experiment, which allows us to work with the probability distribution induced on the resulting set of numbers. Random variables can also be defined for uncountably infinite sample spaces, but they raise technical issues that are unnecessary to address for our purposes. Henceforth, we shall assume that random variables are discrete.

For a random variable $X$ and a real number $x$, we define the event $X = x$ to be $\{s \in S : X(s) = x\}$; thus,

$$\Pr\{X = x\} = \sum_{\{s \in S : X(s) = x\}} \Pr\{s\} \ .$$

The function

$$f(x) = \Pr\{X = x\}$$

is the ***probability density function*** of the random variable $X$. From the probability axioms, $\Pr\{X = x\} \geq 0$ and $\sum_x \Pr\{X = x\} = 1$.

As an example, consider the experiment of rolling a pair of ordinary, 6-sided dice. There are 36 possible elementary events in the sample space. We assume that the probability distribution is uniform, so that each elementary event $s \in S$ is equally likely: $\Pr\{s\} = 1/36$. Define the random variable $X$ to be the *maximum* of the two values showing on the dice. We have $\Pr\{X = 3\} = 5/36$, since $X$ assigns a value of 3 to 5 of the 36 possible elementary events, namely, $(1, 3)$, $(2, 3)$, $(3, 3)$, $(3, 2)$, and $(3, 1)$.

It is common for several random variables to be defined on the same sample space. If $X$ and $Y$ are random variables, the function

$$f(x, y) = \Pr\{X = x \text{ and } Y = y\}$$

is the ***joint probability density function*** of $X$ and $Y$. For a fixed value $y$,

$$\Pr\{Y = y\} = \sum_x \Pr\{X = x \text{ and } Y = y\} \ ,$$

and similarly, for a fixed value $x$,

$$\Pr\{X = x\} = \sum_y \Pr\{X = x \text{ and } Y = y\} \ .$$

Using the definition (C.14) of conditional probability, we have

$$\Pr\{X = x \mid Y = y\} = \frac{\Pr\{X = x \text{ and } Y = y\}}{\Pr\{Y = y\}} \ .$$

We define two random variables $X$ and $Y$ to be ***independent*** if for all $x$ and $y$, the events $X = x$ and $Y = y$ are independent or, equivalently, if for all $x$ and $y$, we have $\Pr\{X = x \text{ and } Y = y\} = \Pr\{X = x\} \Pr\{Y = y\}$.

Given a set of random variables defined over the same sample space, one can define new random variables as sums, products, or other functions of the original variables.

**Expected value of a random variable**

The simplest and most useful summary of the distribution of a random variable is the "average" of the values it takes on. The ***expected value*** (or, synonymously, ***expectation*** or ***mean***) of a discrete random variable $X$ is

$$E[X] = \sum_x x \, \Pr\{X = x\} \, , \tag{C.19}$$

which is well defined if the sum is finite or converges absolutely. Sometimes the expectation of $X$ is denoted by $\mu_X$ or, when the random variable is apparent from context, simply by $\mu$.

Consider a game in which you flip two fair coins. You earn \$3 for each head but lose \$2 for each tail. The expected value of the random variable $X$ representing your earnings is

$$\begin{aligned} E[X] &= 6 \cdot \Pr\{2 \text{ H's}\} + 1 \cdot \Pr\{1 \text{ H}, 1 \text{ T}\} - 4 \cdot \Pr\{2 \text{ T's}\} \\ &= 6(1/4) + 1(1/2) - 4(1/4) \\ &= 1 \, . \end{aligned}$$

The expectation of the sum of two random variables is the sum of their expectations, that is,

$$E[X + Y] = E[X] + E[Y] \, , \tag{C.20}$$

whenever $E[X]$ and $E[Y]$ are defined. We call this property ***linearity of expectation***, and it holds even if $X$ and $Y$ are not independent. It also extends to finite and absolutely convergent summations of expectations. Linearity of expectation is the key property that enables us to perform probabilistic analyses by using indicator random variables (see Section 5.2).

If $X$ is any random variable, any function $g(x)$ defines a new random variable $g(X)$. If the expectation of $g(X)$ is defined, then

$$E[g(X)] = \sum_x g(x) \, \Pr\{X = x\} \, .$$

Letting $g(x) = ax$, we have for any constant $a$,

$$E[aX] = aE[X] \, . \tag{C.21}$$

Consequently, expectations are linear: for any two random variables $X$ and $Y$ and any constant $a$,

$$E[aX + Y] = aE[X] + E[Y] \, . \tag{C.22}$$

When two random variables $X$ and $Y$ are independent and each has a defined expectation,

$$E[XY] = \sum_x \sum_y xy \, \Pr\{X = x \text{ and } Y = y\}$$

$$= \sum_x \sum_y xy \Pr\{X = x\} \Pr\{Y = y\}$$

$$= \left(\sum_x x \Pr\{X = x\}\right)\left(\sum_y y \Pr\{Y = y\}\right)$$

$$= \mathrm{E}[X]\,\mathrm{E}[Y]\ .$$

In general, when $n$ random variables $X_1, X_2, \ldots, X_n$ are mutually independent,

$$\mathrm{E}[X_1 X_2 \cdots X_n] = \mathrm{E}[X_1]\,\mathrm{E}[X_2] \cdots \mathrm{E}[X_n]\ . \tag{C.23}$$

When a random variable $X$ takes on values from the set of natural numbers $\mathbf{N} = \{0, 1, 2, \ldots\}$, there is a nice formula for its expectation:

$$\begin{aligned}
\mathrm{E}[X] &= \sum_{i=0}^{\infty} i\ \Pr\{X = i\} \\
&= \sum_{i=0}^{\infty} i\,(\Pr\{X \geq i\} - \Pr\{X \geq i + 1\}) \\
&= \sum_{i=1}^{\infty} \Pr\{X \geq i\}\ , 
\end{aligned} \tag{C.24}$$

since each term $\Pr\{X \geq i\}$ is added in $i$ times and subtracted out $i - 1$ times (except $\Pr\{X \geq 0\}$, which is added in 0 times and not subtracted out at all).

When we apply a convex function $f(x)$ to a random variable $X$, **Jensen's inequality** gives us

$$\mathrm{E}[f(X)] \geq f(\mathrm{E}[X])\ , \tag{C.25}$$

provided that the expectations exist and are finite. (A function $f(x)$ is **convex** if for all $x$ and $y$ and for all $0 \leq \lambda \leq 1$, we have $f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$.)

### Variance and standard deviation

The expected value of a random variable does not tell us how "spread out" the variable's values are. For example, if we have random variables $X$ and $Y$ for which $\Pr\{X = 1/4\} = \Pr\{X = 3/4\} = 1/2$ and $\Pr\{Y = 0\} = \Pr\{Y = 1\} = 1/2$, then both $\mathrm{E}[X]$ and $\mathrm{E}[Y]$ are $1/2$, yet the actual values taken on by $Y$ are farther from the mean than the actual values taken on by $X$.

The notion of variance mathematically expresses how far from the mean a random variable's values are likely to be. The **variance** of a random variable $X$ with mean $\mathrm{E}[X]$ is

$$
\begin{aligned}
\mathrm{Var}[X] &= \mathrm{E}[(X - \mathrm{E}[X])^2] \\
&= \mathrm{E}[X^2 - 2X\mathrm{E}[X] + \mathrm{E}^2[X]] \\
&= \mathrm{E}[X^2] - 2\mathrm{E}[X\mathrm{E}[X]] + \mathrm{E}^2[X] \\
&= \mathrm{E}[X^2] - 2\mathrm{E}^2[X] + \mathrm{E}^2[X] \\
&= \mathrm{E}[X^2] - \mathrm{E}^2[X] \ . \tag{C.26}
\end{aligned}
$$

The justification for the equalities $\mathrm{E}[\mathrm{E}^2[X]] = \mathrm{E}^2[X]$ and $\mathrm{E}[X\mathrm{E}[X]] = \mathrm{E}^2[X]$ is that $\mathrm{E}[X]$ is not a random variable but simply a real number, which means that equation (C.21) applies (with $a = \mathrm{E}[X]$). Equation (C.26) can be rewritten to obtain an expression for the expectation of the square of a random variable:

$$
\mathrm{E}[X^2] = \mathrm{Var}[X] + \mathrm{E}^2[X] \ . \tag{C.27}
$$

The variance of a random variable $X$ and the variance of $aX$ are related (see Exercise C.3-10):

$$
\mathrm{Var}[aX] = a^2\mathrm{Var}[X] \ .
$$

When $X$ and $Y$ are independent random variables,

$$
\mathrm{Var}[X + Y] = \mathrm{Var}[X] + \mathrm{Var}[Y] \ .
$$

In general, if $n$ random variables $X_1, X_2, \ldots, X_n$ are pairwise independent, then

$$
\mathrm{Var}\left[\sum_{i=1}^{n} X_i\right] = \sum_{i=1}^{n} \mathrm{Var}[X_i] \ . \tag{C.28}
$$

The ***standard deviation*** of a random variable $X$ is the positive square root of the variance of $X$. The standard deviation of a random variable $X$ is sometimes denoted $\sigma_X$ or simply $\sigma$ when the random variable $X$ is understood from context. With this notation, the variance of $X$ is denoted $\sigma^2$.

**Exercises**

***C.3-1***
Two ordinary, 6-sided dice are rolled. What is the expectation of the sum of the two values showing? What is the expectation of the maximum of the two values showing?

***C.3-2***
An array $A[1 .. n]$ contains $n$ distinct numbers that are randomly ordered, with each permutation of the $n$ numbers being equally likely. What is the expectation of the index of the maximum element in the array? What is the expectation of the index of the minimum element in the array?

### C.3-3
A carnival game consists of three dice in a cage. A player can bet a dollar on any of the numbers 1 through 6. The cage is shaken, and the payoff is as follows. If the player's number doesn't appear on any of the dice, he loses his dollar. Otherwise, if his number appears on exactly $k$ of the three dice, for $k = 1, 2, 3$, he keeps his dollar and wins $k$ more dollars. What is his expected gain from playing the carnival game once?

### C.3-4
Argue that if $X$ and $Y$ are nonnegative random variables, then

$$\mathrm{E}\left[\max(X, Y)\right] \leq \mathrm{E}\left[X\right] + \mathrm{E}\left[Y\right] .$$

### C.3-5  ★
Let $X$ and $Y$ be independent random variables. Prove that $f(X)$ and $g(Y)$ are independent for any choice of functions $f$ and $g$.

### C.3-6  ★
Let $X$ be a nonnegative random variable, and suppose that $\mathrm{E}\left[X\right]$ is well defined. Prove ***Markov's inequality***:

$$\Pr\{X \geq t\} \leq \mathrm{E}\left[X\right]/t \tag{C.29}$$

for all $t > 0$.

### C.3-7  ★
Let $S$ be a sample space, and let $X$ and $X'$ be random variables such that $X(s) \geq X'(s)$ for all $s \in S$. Prove that for any real constant $t$,

$$\Pr\{X \geq t\} \geq \Pr\{X' \geq t\} .$$

### C.3-8
Which is larger: the expectation of the square of a random variable, or the square of its expectation?

### C.3-9
Show that for any random variable $X$ that takes on only the values 0 and 1, we have $\mathrm{Var}\left[X\right] = \mathrm{E}\left[X\right]\mathrm{E}\left[1 - X\right]$.

### C.3-10
Prove that $\mathrm{Var}\left[aX\right] = a^2\mathrm{Var}\left[X\right]$ from the definition (C.26) of variance.

## C.4    The geometric and binomial distributions

A coin flip is an instance of a **Bernoulli trial**, which is defined as an experiment with only two possible outcomes: **success**, which occurs with probability $p$, and **failure**, which occurs with probability $q = 1 - p$. When we speak of **Bernoulli trials** collectively, we mean that the trials are mutually independent and, unless we specifically say otherwise, that each has the same probability $p$ for success. Two important distributions arise from Bernoulli trials: the geometric distribution and the binomial distribution.

### The geometric distribution

Suppose we have a sequence of Bernoulli trials, each with a probability $p$ of success and a probability $q = 1 - p$ of failure. How many trials occur before we obtain a success? Let the random variable $X$ be the number of trials needed to obtain a success. Then $X$ has values in the range $\{1, 2, \ldots\}$, and for $k \geq 1$,

$$\Pr\{X = k\} = q^{k-1}p \,, \tag{C.30}$$

since we have $k - 1$ failures before the one success. A probability distribution satisfying equation (C.30) is said to be a **geometric distribution**. Figure C.1 illustrates such a distribution.

Assuming that $q < 1$, the expectation of a geometric distribution can be calculated using identity (A.8):

$$\begin{aligned}
\mathrm{E}[X] &= \sum_{k=1}^{\infty} kq^{k-1}p \\
&= \frac{p}{q} \sum_{k=0}^{\infty} kq^{k} \\
&= \frac{p}{q} \cdot \frac{q}{(1-q)^2} \\
&= 1/p \,.
\end{aligned} \tag{C.31}$$

Thus, on average, it takes $1/p$ trials before we obtain a success, an intuitive result. The variance, which can be calculated similarly, but using Exercise A.1-3, is

$$\mathrm{Var}[X] = q/p^2 \,. \tag{C.32}$$

As an example, suppose we repeatedly roll two dice until we obtain either a seven or an eleven. Of the 36 possible outcomes, 6 yield a seven and 2 yield an eleven. Thus, the probability of success is $p = 8/36 = 2/9$, and we must roll $1/p = 9/2 = 4.5$ times on average to obtain a seven or eleven.
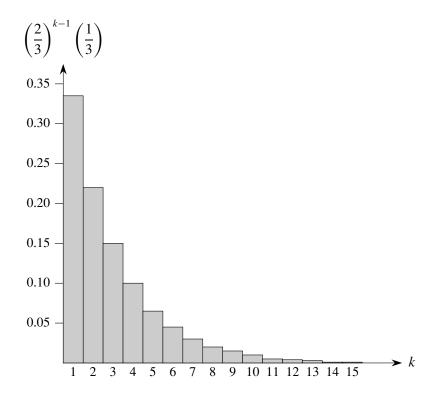
$$\left(\frac{2}{3}\right)^{k-1}\left(\frac{1}{3}\right)$$



**Figure C.1** A geometric distribution with probability $p = 1/3$ of success and a probability $q = 1 - p$ of failure. The expectation of the distribution is $1/p = 3$.

### The binomial distribution

How many successes occur during $n$ Bernoulli trials, where a success occurs with probability $p$ and a failure with probability $q = 1 - p$? Define the random variable $X$ to be the number of successes in $n$ trials. Then $X$ has values in the range $\{0, 1, \ldots, n\}$, and for $k = 0, \ldots, n$,

$$\Pr\{X = k\} = \binom{n}{k}p^k q^{n-k} , \tag{C.33}$$

since there are $\binom{n}{k}$ ways to pick which $k$ of the $n$ trials are successes, and the probability that each occurs is $p^k q^{n-k}$. A probability distribution satisfying equation (C.33) is said to be a ***binomial distribution***. For convenience, we define the family of binomial distributions using the notation

$$b(k; n, p) = \binom{n}{k}p^k (1 - p)^{n-k} . \tag{C.34}$$

Figure C.2 illustrates a binomial distribution. The name "binomial" comes from the fact that (C.33) is the $k$th term of the expansion of $(p + q)^n$. Consequently, since $p + q = 1$,
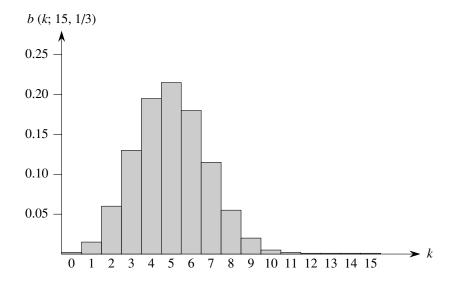
$b\,(k;\,15,\,1/3)$



**Figure C.2**   The binomial distribution $b(k; 15, 1/3)$ resulting from $n = 15$ Bernoulli trials, each with probability $p = 1/3$ of success. The expectation of the distribution is $np = 5$.

$$\sum_{k=0}^{n} b(k; n, p) = 1 \,, \tag{C.35}$$

as is required by axiom 2 of the probability axioms.

   We can compute the expectation of a random variable having a binomial distribution from equations (C.8) and (C.35). Let $X$ be a random variable that follows the binomial distribution $b(k; n, p)$, and let $q = 1 - p$. By the definition of expectation, we have

$$
\begin{aligned}
\mathrm{E}\,[X] &= \sum_{k=0}^{n} k \, \mathrm{Pr}\,\{X = k\} \\
&= \sum_{k=0}^{n} k \, b(k; n, p) \\
&= \sum_{k=1}^{n} k \binom{n}{k} p^k q^{n-k} \\
&= np \sum_{k=1}^{n} \binom{n-1}{k-1} p^{k-1} q^{n-k} \\
&= np \sum_{k=0}^{n-1} \binom{n-1}{k} p^k q^{(n-1)-k}
\end{aligned}
$$

$$= np \sum_{k=0}^{n-1} b(k; n-1, p)$$
$$= np . \tag{C.36}$$

By using the linearity of expectation, we can obtain the same result with substantially less algebra. Let $X_i$ be the random variable describing the number of successes in the $i$th trial. Then $\mathrm{E}[X_i] = p \cdot 1 + q \cdot 0 = p$, and by linearity of expectation (equation (C.20)), the expected number of successes for $n$ trials is

$$
\begin{aligned}
\mathrm{E}[X] &= \mathrm{E}\left[\sum_{i=1}^{n} X_i\right] \\
&= \sum_{i=1}^{n} \mathrm{E}[X_i] \\
&= \sum_{i=1}^{n} p \\
&= np .
\end{aligned}
\tag{C.37}
$$

The same approach can be used to calculate the variance of the distribution. Using equation (C.26), we have $\mathrm{Var}[X_i] = \mathrm{E}[X_i^2] - \mathrm{E}^2[X_i]$. Since $X_i$ only takes on the values 0 and 1, we have $\mathrm{E}[X_i^2] = \mathrm{E}[X_i] = p$, and hence

$$\mathrm{Var}[X_i] = p - p^2 = pq . \tag{C.38}$$

To compute the variance of $X$, we take advantage of the independence of the $n$ trials; thus, by equation (C.28),

$$
\begin{aligned}
\mathrm{Var}[X] &= \mathrm{Var}\left[\sum_{i=1}^{n} X_i\right] \\
&= \sum_{i=1}^{n} \mathrm{Var}[X_i] \\
&= \sum_{i=1}^{n} pq \\
&= npq .
\end{aligned}
\tag{C.39}
$$

As can be seen from Figure C.2, the binomial distribution $b(k; n, p)$ increases as $k$ runs from 0 to $n$ until it reaches the mean $np$, and then it decreases. We can prove that the distribution always behaves in this manner by looking at the ratio of successive terms:

$$
\begin{aligned}
\frac{b(k; n, p)}{b(k - 1; n, p)} &= \frac{\binom{n}{k} p^k q^{n-k}}{\binom{n}{k-1} p^{k-1} q^{n-k+1}} \\
&= \frac{n!(k-1)!(n-k+1)!\, p}{k!(n-k)!\, n!\, q} \\
&= \frac{(n-k+1)p}{kq} \\
&= 1 + \frac{(n+1)p - k}{kq} .
\end{aligned}
\tag{C.40}
$$

This ratio is greater than 1 precisely when $(n+1)p - k$ is positive. Consequently, $b(k; n, p) > b(k - 1; n, p)$ for $k < (n+1)p$ (the distribution increases), and $b(k; n, p) < b(k - 1; n, p)$ for $k > (n+1)p$ (the distribution decreases). If $k = (n+1)p$ is an integer, then $b(k; n, p) = b(k - 1; n, p)$, so the distribution has two maxima: at $k = (n+1)p$ and at $k - 1 = (n+1)p - 1 = np - q$. Otherwise, it attains a maximum at the unique integer $k$ that lies in the range $np - q < k < (n+1)p$.

The following lemma provides an upper bound on the binomial distribution.

***Lemma C.1***
Let $n \geq 0$, let $0 < p < 1$, let $q = 1 - p$, and let $0 \leq k \leq n$. Then

$$
b(k; n, p) \leq \left(\frac{np}{k}\right)^k \left(\frac{nq}{n - k}\right)^{n-k} .
$$

***Proof***   Using equation (C.6), we have

$$
\begin{aligned}
b(k; n, p) &= \binom{n}{k} p^k q^{n-k} \\
&\leq \left(\frac{n}{k}\right)^k \left(\frac{n}{n - k}\right)^{n-k} p^k q^{n-k} \\
&= \left(\frac{np}{k}\right)^k \left(\frac{nq}{n - k}\right)^{n-k} .
\end{aligned}
$$
■

### Exercises

***C.4-1***
Verify axiom 2 of the probability axioms for the geometric distribution.

***C.4-2***
How many times on average must we flip 6 fair coins before we obtain 3 heads and 3 tails?

***C.4-3***
Show that $b(k; n, p) = b(n - k; n, q)$, where $q = 1 - p$.

### C.4-4

Show that value of the maximum of the binomial distribution $b(k; n, p)$ is approximately $1/\sqrt{2\pi npq}$, where $q = 1 - p$.

### C.4-5   ★

Show that the probability of no successes in $n$ Bernoulli trials, each with probability $p = 1/n$, is approximately $1/e$. Show that the probability of exactly one success is also approximately $1/e$.

### C.4-6   ★

Professor Rosencrantz flips a fair coin $n$ times, and so does Professor Guildenstern. Show that the probability that they get the same number of heads is $\binom{2n}{n}/4^n$. (*Hint:* For Professor Rosencrantz, call a head a success; for Professor Guildenstern, call a tail a success.) Use your argument to verify the identity

$$\sum_{k=0}^{n} \binom{n}{k}^2 = \binom{2n}{n} .$$

### C.4-7   ★

Show that for $0 \le k \le n$,

$$b(k; n, 1/2) \le 2^{n\,H(k/n)-n} ,$$

where $H(x)$ is the entropy function (C.7).

### C.4-8   ★

Consider $n$ Bernoulli trials, where for $i = 1, 2, \ldots, n$, the $i$th trial has probability $p_i$ of success, and let $X$ be the random variable denoting the total number of successes. Let $p \ge p_i$ for all $i = 1, 2, \ldots, n$. Prove that for $1 \le k \le n$,

$$\Pr\{X < k\} \le \sum_{i=0}^{k-1} b(i; n, p) .$$

### C.4-9   ★

Let $X$ be the random variable for the total number of successes in a set $A$ of $n$ Bernoulli trials, where the $i$th trial has a probability $p_i$ of success, and let $X'$ be the random variable for the total number of successes in a second set $A'$ of $n$ Bernoulli trials, where the $i$th trial has a probability $p'_i \ge p_i$ of success. Prove that for $0 \le k \le n$,

$$\Pr\{X' \ge k\} \ge \Pr\{X \ge k\} .$$

(*Hint:* Show how to obtain the Bernoulli trials in $A'$ by an experiment involving the trials of $A$, and use the result of Exercise C.3-7.)

## ★  C.5    The tails of the binomial distribution

The probability of having at least, or at most, $k$ successes in $n$ Bernoulli trials, each with probability $p$ of success, is often of more interest than the probability of having exactly $k$ successes. In this section, we investigate the ***tails*** of the binomial distribution: the two regions of the distribution $b(k; n, p)$ that are far from the mean $np$. We shall prove several important bounds on (the sum of all terms in) a tail.

We first provide a bound on the right tail of the distribution $b(k; n, p)$. Bounds on the left tail can be determined by inverting the roles of successes and failures.

***Theorem C.2***
Consider a sequence of $n$ Bernoulli trials, where success occurs with probability $p$. Let $X$ be the random variable denoting the total number of successes. Then for $0 \le k \le n$, the probability of at least $k$ successes is

$$\Pr\{X \ge k\} = \sum_{i=k}^{n} b(i; n, p)$$

$$\le \binom{n}{k} p^k .$$

***Proof***    For $S \subseteq \{1, 2, \ldots, n\}$, we let $A_S$ denote the event that the $i$th trial is a success for every $i \in S$. Clearly $\Pr\{A_S\} = p^k$ if $|S| = k$. We have

$$\Pr\{X \ge k\} = \Pr\{\text{there exists } S \subseteq \{1, 2, \ldots, n\} : |S| = k \text{ and } A_S\}$$

$$= \Pr\left\{\bigcup_{S \subseteq \{1,2,\ldots,n\}:|S|=k} A_S\right\}$$

$$\le \sum_{S \subseteq \{1,2,\ldots,n\}:|S|=k} \Pr\{A_S\}$$

$$= \binom{n}{k} p^k ,$$

where the inequality follows from Boole's inequality (C.18).    ∎

The following corollary restates the theorem for the left tail of the binomial distribution. In general, we shall leave it to the reader to adapt the proofs from one tail to the other.

***Corollary C.3***
Consider a sequence of $n$ Bernoulli trials, where success occurs with probability $p$. If $X$ is the random variable denoting the total number of successes, then for

$0 \le k \le n$, the probability of at most $k$ successes is

$$
\begin{aligned}
\Pr\{X \le k\} &= \sum_{i=0}^{k} b(i; n, p) \\
&\le \binom{n}{n-k}(1-p)^{n-k} \\
&= \binom{n}{k}(1-p)^{n-k} .
\end{aligned}
$$
■

Our next bound concerns the left tail of the binomial distribution. Its corollary shows that, far from the mean, the left tail diminishes exponentially.

***Theorem C.4***
Consider a sequence of $n$ Bernoulli trials, where success occurs with probability $p$ and failure with probability $q = 1 - p$. Let $X$ be the random variable denoting the total number of successes. Then for $0 < k < np$, the probability of fewer than $k$ successes is

$$
\begin{aligned}
\Pr\{X < k\} &= \sum_{i=0}^{k-1} b(i; n, p) \\
&< \frac{kq}{np-k} b(k; n, p) .
\end{aligned}
$$

***Proof***   We bound the series $\sum_{i=0}^{k-1} b(i; n, p)$ by a geometric series using the technique from Section A.2, page 1064. For $i = 1, 2, \ldots, k$, we have from equation (C.40),

$$
\begin{aligned}
\frac{b(i-1; n, p)}{b(i; n, p)} &= \frac{iq}{(n-i+1)p} \\
&< \frac{iq}{(n-i)p} \\
&\le \frac{kq}{(n-k)p} .
\end{aligned}
$$

If we let

$$
\begin{aligned}
x &= \frac{kq}{(n-k)p} \\
&< \frac{kq}{(n-np)p} \\
&= \frac{kq}{nqp}
\end{aligned}
$$

$$= \frac{k}{np}$$
$$< 1 ,$$

it follows that

$$b(i - 1; n, p) < x \, b(i; n, p)$$

for $0 < i \le k$. Iteratively applying this inequality $k - i$ times, we obtain

$$b(i; n, p) < x^{k-i} \, b(k; n, p)$$

for $0 \le i < k$, and hence

$$
\sum_{i=0}^{k-1} b(i; n, p) \;\; < \;\; \sum_{i=0}^{k-1} x^{k-i} b(k; n, p)
$$
$$
< \;\; b(k; n, p) \sum_{i=1}^{\infty} x^i
$$
$$
= \;\; \frac{x}{1 - x} b(k; n, p)
$$
$$
= \;\; \frac{kq}{np - k} b(k; n, p) \; . \qquad \blacksquare
$$

### Corollary C.5

Consider a sequence of $n$ Bernoulli trials, where success occurs with probability $p$ and failure with probability $q = 1 - p$. Then for $0 < k \le np/2$, the probability of fewer than $k$ successes is less than one half of the probability of fewer than $k + 1$ successes.

***Proof***   Because $k \le np/2$, we have

$$
\frac{kq}{np - k} \;\; \le \;\; \frac{(np/2)q}{np - (np/2)}
$$
$$
= \;\; \frac{(np/2)q}{np/2}
$$
$$
\le \;\; 1 ,
$$

since $q \le 1$. Letting $X$ be the random variable denoting the number of successes, Theorem C.4 implies that the probability of fewer than $k$ successes is

$$
\Pr\{X < k\} = \sum_{i=0}^{k-1} b(i; n, p) < b(k; n, p) \; .
$$

Thus we have

$$
\begin{aligned}
\frac{\Pr\{X < k\}}{\Pr\{X < k+1\}} &= \frac{\sum_{i=0}^{k-1} b(i; n, p)}{\sum_{i=0}^{k} b(i; n, p)} \\
&= \frac{\sum_{i=0}^{k-1} b(i; n, p)}{\sum_{i=0}^{k-1} b(i; n, p) + b(k; n, p)} \\
&< 1/2 \,,
\end{aligned}
$$

since $\sum_{i=0}^{k-1} b(i; n, p) < b(k; n, p)$. ∎

Bounds on the right tail can be determined similarly. Their proofs are left as Exercise C.5-2.

### Corollary C.6
Consider a sequence of $n$ Bernoulli trials, where success occurs with probability $p$. Let $X$ be the random variable denoting the total number of successes. Then for $np < k < n$, the probability of more than $k$ successes is

$$
\begin{aligned}
\Pr\{X > k\} &= \sum_{i=k+1}^{n} b(i; n, p) \\
&< \frac{(n-k)p}{k - np} b(k; n, p) \,.
\end{aligned}
$$
∎

### Corollary C.7
Consider a sequence of $n$ Bernoulli trials, where success occurs with probability $p$ and failure with probability $q = 1 - p$. Then for $(np + n)/2 < k < n$, the probability of more than $k$ successes is less than one half of the probability of more than $k - 1$ successes. ∎

The next theorem considers $n$ Bernoulli trials, each with a probability $p_i$ of success, for $i = 1, 2, \ldots, n$. As the subsequent corollary shows, we can use the theorem to provide a bound on the right tail of the binomial distribution by setting $p_i = p$ for each trial.

### Theorem C.8
Consider a sequence of $n$ Bernoulli trials, where in the $i$th trial, for $i = 1, 2, \ldots, n$, success occurs with probability $p_i$ and failure occurs with probability $q_i = 1 - p_i$. Let $X$ be the random variable describing the total number of successes, and let $\mu = \mathrm{E}[X]$. Then for $r > \mu$,

$$
\Pr\{X - \mu \geq r\} \leq \left(\frac{\mu e}{r}\right)^r \,.
$$

***Proof***   Since for any $\alpha > 0$, the function $e^{\alpha x}$ is strictly increasing in $x$,

$$\Pr\{X - \mu \geq r\} = \Pr\{e^{\alpha(X-\mu)} \geq e^{\alpha r}\} \ , \tag{C.41}$$

where $\alpha$ will be determined later. Using Markov's inequality (C.29), we obtain

$$\Pr\{e^{\alpha(X-\mu)} \geq e^{\alpha r}\} \leq \mathrm{E}[e^{\alpha(X-\mu)}]e^{-\alpha r} \ . \tag{C.42}$$

The bulk of the proof consists of bounding $\mathrm{E}[e^{\alpha(X-\mu)}]$ and substituting a suitable value for $\alpha$ in inequality (C.42). First, we evaluate $\mathrm{E}[e^{\alpha(X-\mu)}]$. Using the notation of Section 5.2, let $X_i = \mathrm{I}\{\text{the } i\text{th Bernoulli trial is a success}\}$ for $i = 1, 2, \ldots, n$; that is, $X_i$ is the random variable that is 1 if the $i$th Bernoulli trial is a success and 0 if it is a failure. Thus,

$$X = \sum_{i=1}^{n} X_i \ ,$$

and by linearity of expectation,

$$\mu = \mathrm{E}[X] = \mathrm{E}\left[\sum_{i=1}^{n} X_i\right] = \sum_{i=1}^{n} \mathrm{E}[X_i] = \sum_{i=1}^{n} p_i \ ,$$

which implies

$$X - \mu = \sum_{i=1}^{n}(X_i - p_i) \ .$$

To evaluate $\mathrm{E}[e^{\alpha(X-\mu)}]$, we substitute for $X - \mu$, obtaining

$$
\begin{aligned}
\mathrm{E}[e^{\alpha(X-\mu)}] &= \mathrm{E}\left[e^{\alpha \sum_{i=1}^{n}(X_i - p_i)}\right] \\
&= \mathrm{E}\left[\prod_{i=1}^{n} e^{\alpha(X_i - p_i)}\right] \\
&= \prod_{i=1}^{n} \mathrm{E}[e^{\alpha(X_i - p_i)}] \ ,
\end{aligned}
$$

which follows from (C.23), since the mutual independence of the random variables $X_i$ implies the mutual independence of the random variables $e^{\alpha(X_i - p_i)}$ (see Exercise C.3-5). By the definition of expectation,

$$
\begin{aligned}
\mathrm{E}[e^{\alpha(X_i - p_i)}] &= e^{\alpha(1-p_i)}p_i + e^{\alpha(0-p_i)}q_i \\
&= p_i e^{\alpha q_i} + q_i e^{-\alpha p_i} \\
&\leq p_i e^{\alpha} + 1 \\
&\leq \exp(p_i e^{\alpha}) \ ,
\end{aligned}
\tag{C.43}
$$

where $\exp(x)$ denotes the exponential function: $\exp(x) = e^x$. (Inequality (C.43) follows from the inequalities $\alpha > 0$, $q_i \leq 1$, $e^{\alpha q_i} \leq e^\alpha$, and $e^{-\alpha p_i} \leq 1$, and the last line follows from inequality (3.11)). Consequently,

$$
\begin{aligned}
\mathrm{E}\left[e^{\alpha(X-\mu)}\right] &= \prod_{i=1}^{n} \mathrm{E}\left[e^{\alpha(X_i-p_i)}\right] \\
&\leq \prod_{i=1}^{n} \exp(p_i e^\alpha) \\
&= \exp\left(\sum_{i=1}^{n} p_i e^\alpha\right) \\
&= \exp(\mu e^\alpha) , \quad\quad\quad\quad\quad\quad\quad\quad \text{(C.44)}
\end{aligned}
$$

since $\mu = \sum_{i=1}^{n} p_i$. Therefore, from equation (C.41) and inequalities (C.42) and (C.44), it follows that

$$
\Pr\{X - \mu \geq r\} \leq \exp(\mu e^\alpha - \alpha r) . \quad\quad\quad\quad\quad\quad \text{(C.45)}
$$

Choosing $\alpha = \ln(r/\mu)$ (see Exercise C.5-7), we obtain

$$
\begin{aligned}
\Pr\{X - \mu \geq r\} &\leq \exp(\mu e^{\ln(r/\mu)} - r \ln(r/\mu)) \\
&= \exp(r - r \ln(r/\mu)) \\
&= \frac{e^r}{(r/\mu)^r} \\
&= \left(\frac{\mu e}{r}\right)^r . \quad\quad\quad\quad\quad\quad\quad\quad\blacksquare
\end{aligned}
$$

When applied to Bernoulli trials in which each trial has the same probability of success, Theorem C.8 yields the following corollary bounding the right tail of a binomial distribution.

### Corollary C.9
Consider a sequence of $n$ Bernoulli trials, where in each trial success occurs with probability $p$ and failure occurs with probability $q = 1 - p$. Then for $r > np$,

$$
\begin{aligned}
\Pr\{X - np \geq r\} &= \sum_{k=\lceil np+r \rceil}^{n} b(k; n, p) \\
&\leq \left(\frac{npe}{r}\right)^r .
\end{aligned}
$$

***Proof*** By equation (C.36), we have $\mu = \mathrm{E}[X] = np$. $\quad\quad\quad\quad\quad\quad\blacksquare$

**Exercises**

*C.5-1*  ⋆
Which is less likely: obtaining no heads when you flip a fair coin $n$ times, or obtaining fewer than $n$ heads when you flip the coin $4n$ times?

*C.5-2*  ⋆
Prove Corollaries C.6 and C.7.

*C.5-3*  ⋆
Show that

$$\sum_{i=0}^{k-1} \binom{n}{i} a^i < (a+1)^n \frac{k}{na - k(a+1)} b(k; n, a/(a+1))$$

for all $a > 0$ and all $k$ such that $0 < k < n$.

*C.5-4*  ⋆
Prove that if $0 < k < np$, where $0 < p < 1$ and $q = 1 - p$, then

$$\sum_{i=0}^{k-1} p^i q^{n-i} < \frac{kq}{np - k} \left(\frac{np}{k}\right)^k \left(\frac{nq}{n-k}\right)^{n-k} .$$

*C.5-5*  ⋆
Show that the conditions of Theorem C.8 imply that

$$\Pr\{\mu - X \geq r\} \leq \left(\frac{(n - \mu)e}{r}\right)^r .$$

Similarly, show that the conditions of Corollary C.9 imply that

$$\Pr\{np - X \geq r\} \leq \left(\frac{nqe}{r}\right)^r .$$

*C.5-6*  ⋆
Consider a sequence of $n$ Bernoulli trials, where in the $i$th trial, for $i = 1, 2, \ldots, n$, success occurs with probability $p_i$ and failure occurs with probability $q_i = 1 - p_i$. Let $X$ be the random variable describing the total number of successes, and let $\mu = \mathrm{E}[X]$. Show that for $r \geq 0$,

$$\Pr\{X - \mu \geq r\} \leq e^{-r^2/2n} .$$

(*Hint:* Prove that $p_i e^{\alpha q_i} + q_i e^{-\alpha p_i} \leq e^{\alpha^2/2}$. Then follow the outline of the proof of Theorem C.8, using this inequality in place of inequality (C.43).)

### C.5-7 ★
Show that the right-hand side of inequality (C.45) is minimized by choosing $\alpha = \ln(r/\mu)$.

## Problems

### C-1 Balls and bins
In this problem, we investigate the effect of various assumptions on the number of ways of placing $n$ balls into $b$ distinct bins.

**a.** Suppose that the $n$ balls are distinct and that their order within a bin does not matter. Argue that the number of ways of placing the balls in the bins is $b^n$.

**b.** Suppose that the balls are distinct and that the balls in each bin are ordered. Prove that there are exactly $(b + n - 1)!/(b - 1)!$ ways to place the balls in the bins. (*Hint:* Consider the number of ways of arranging $n$ distinct balls and $b - 1$ indistinguishable sticks in a row.)

**c.** Suppose that the balls are identical, and hence their order within a bin does not matter. Show that the number of ways of placing the balls in the bins is $\binom{b+n-1}{n}$. (*Hint:* Of the arrangements in part (b), how many are repeated if the balls are made identical?)

**d.** Suppose that the balls are identical and that no bin may contain more than one ball. Show that the number of ways of placing the balls is $\binom{b}{n}$.

**e.** Suppose that the balls are identical and that no bin may be left empty. Show that the number of ways of placing the balls is $\binom{n-1}{b-1}$.

## Chapter notes

The first general methods for solving probability problems were discussed in a famous correspondence between B. Pascal and P. de Fermat, which began in 1654, and in a book by C. Huygens in 1657. Rigorous probability theory began with the work of J. Bernoulli in 1713 and A. De Moivre in 1730. Further developments of the theory were provided by P. S. de Laplace, S.-D. Poisson, and C. F. Gauss.

Sums of random variables were originally studied by P. L. Chebyshev and A. A. Markov. Probability theory was axiomatized by A. N. Kolmogorov in 1933.

Bounds on the tails of distributions were provided by Chernoff [59] and Hoeffding [150]. Seminal work in random combinatorial structures was done by P. Erdös.

Knuth [182] and Liu [205] are good references for elementary combinatorics and counting. Standard textbooks such as Billingsley [42], Chung [60], Drake [80], Feller [88], and Rozanov [263] offer comprehensive introductions to probability.