
C Counting and Probability

This appendix reviews elementary combinatorics and probability theory. If you have a good background in these areas, you may want to skim the beginning of this appendix lightly and concentrate on the later sections. Most of this book's 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 form 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 continues 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 all the choices. For example, you might ask, “How many different n -bit numbers are there?” or “How many orderings of n distinct elements are there?” This section reviews the elements of counting theory. Since some of the material assumes a basic understanding of sets, you might wish to start by reviewing the material in Section B.1.

Rules of sum and product

We can sometimes express a set of items that we wish to count 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 one 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) on page 1156. For example, if each position on a car's license plate is a letter or a digit, then the number of possibilities for each position is $26 + 10 = 36$, 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) on page 1157. For example, if an ice-cream parlor offers 28 flavors of ice cream and four 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 eight binary strings of length 3:

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

(Here we use the shorthand of omitting the angle brackets when denoting a sequence.) 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.

We can view a k -string over a set S as an element of the Cartesian product S^k of k -tuples, which means that 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, there are n ways to pick the first element; for each of these choices, there are n ways to pick the second element; and so forth k times. This construction leads to the k -fold product $\underbrace{n \cdot n \cdots n}_{n \text{ times}} = n^k$ 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 $S = \{a, b, c\}$, then S has 6 permutations:

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

(Again, we use the shorthand of omitting the angle brackets when denoting a sequence.) There are $n!$ permutations of a set of n elements, since there are n ways to choose the first element of the sequence, $n - 1$ ways for the second element, $n - 2$ ways for the third, 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 an n -permutation of an n -set.) Here are the 2-permutations of the set $\{a, b, c, d\}$:

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

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

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

since there are n ways to choose the first element, $n-1$ ways to choose the second element, and so on, until k elements are chosen, with the last element chosen from the remaining $n-k+1$ elements. For the above example, with $n=4$ and $k=2$, the formula (C.1) evaluates to $4!/2! = 12$, matching the number of 2-permutations listed.

Combinations

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

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

(Here we use the shorthand of omitting the braces around each subset.) To construct a k -combination of an n -set, choose k distinct (different) elements from the n -set. The order of selecting the elements does not matter.

We can express the number of k -combinations of an n -set in terms of the number of k -permutations of an n -set. Every k -combination has 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

$$\frac{n!}{k!(n-k)!} . \quad (\text{C.2})$$

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

The notation $\binom{n}{k}$ (read “ n choose k ”) denotes the number of k -combinations of an n -set. Equation (C.2) gives

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

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

$$\binom{n}{k} = \binom{n}{n-k}. \quad (\text{C.3})$$

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

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

where $n \in \mathbb{N}$ and $x, y \in \mathbb{R}$. The right-hand side of equation (C.4) is called the *binomial expansion* of the left-hand side. A special case of the binomial theorem 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 1s they contain: $\binom{n}{k}$ binary n -strings contain exactly k 1s, since there are $\binom{n}{k}$ ways to choose k out of the n positions in which to place the 1s.

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 *outcomes* or *elementary events*. Think of each outcome as a possible result of an experiment. For the experiment of flipping two distinguishable coins, with each individual flip resulting in a head (H) or a tail (T), you can view the sample space S as consisting of the set of all possible 2-strings over $\{H, T\}$:

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

An *event* is a subset¹ 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 $A \cap B = \emptyset$. An outcome s also defines the event $\{s\}$, which we sometimes write as just s . By definition, all outcomes 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 satisfying the following *probability axioms*:

1. $\Pr\{A\} \geq 0$ for any event A .
2. $\Pr\{S\} = 1$.
3. $\Pr\{A \cup B\} = \Pr\{A\} + \Pr\{B\}$ for any two mutually exclusive events A and B .
More generally, for any sequence of events A_1, A_2, \dots (finite or countably infinite) 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 . Axiom 2 is simply 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

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 for what subsets are events is that the set of events of a sample space must 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 see in this book are over finite or countable sample spaces, and we generally consider all subsets of a sample space to be events. A notable exception is the continuous uniform probability distribution, which we'll see shortly.

$\Pr\{A\} \leq \Pr\{B\}$. 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 ,

$$\Pr\{A \cup B\} = \Pr\{A\} + \Pr\{B\} - \Pr\{A \cap B\} \quad (\text{C.14})$$

$$\leq \Pr\{A\} + \Pr\{B\} . \quad (\text{C.15})$$

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

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

Another way to obtain the same result is to observe that 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 outcomes, specifically those in A , are mutually exclusive. If S is finite and every outcome $s \in S$ 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, $1/2$. Flipping the coin n times gives the uniform probability distribution defined on the sample space $S = \{\text{H}, \text{T}\}^n$, a set of size 2^n . We can represent each outcome in S as a string of length n over $\{\text{H}, \text{T}\}$, with each string occurring with probability $1/2^n$. The event $A = \{\text{exactly } k \text{ heads and exactly } n - k \text{ tails occur}\}$ is a subset of S of size $|A| = \binom{n}{k}$, since $\binom{n}{k}$ strings of length n over $\{\text{H}, \text{T}\}$ 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 $a < b$. The intuition is that each point in the interval $[a, b]$ should be “equally likely.” Because there are an uncountable number

of points, however, if all points had the same finite, positive probability, axioms 2 and 3 would not be simultaneously satisfied. For this reason, we'd 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

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

Letting $c = d$ gives that the probability of a single point is 0. Removing the endpoints $[c, c]$ and $[d, d]$ of an interval $[c, d]$ results in the open interval (c, d) . Since $[c, d] = [c, c] \cup (c, d) \cup [d, d]$, axiom 3 gives $\Pr\{[c, d]\} = \Pr\{(c, d)\}$. Generally, the set of events for the continuous uniform probability distribution contains any subset of the sample space $[a, b]$ that can be obtained by a finite or countable union of open and closed intervals, as well as certain more complicated sets.

Conditional probability and independence

Sometimes you 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 outcomes are equally likely, and so you infer that each occurs with probability $1/3$. Since only one of these outcomes shows two heads, the answer 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

$$\Pr\{A \mid B\} = \frac{\Pr\{A \cap B\}}{\Pr\{B\}} \quad (\text{C.16})$$

whenever $\Pr\{B\} \neq 0$. (Read " $\Pr\{A \mid B\}$ " as "the probability of A given B .") The idea behind equation (C.16) is that 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. Because 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

$$\Pr\{A \cap B\} = \Pr\{A\} \Pr\{B\} , \quad (\text{C.17})$$

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

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

For example, suppose that you flip two fair coins 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 you 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, \dots, 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 $A_{i_1}, A_{i_2}, \dots, A_{i_k}$ of the collection, where $2 \leq k \leq n$ and $1 \leq i_1 < i_2 < \dots < i_k \leq n$, satisfies

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

For example, suppose that you 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. Then,

$$\Pr\{A_1\} = 1/2 ,$$

$$\Pr\{A_2\} = 1/2 ,$$

$$\Pr\{A_3\} = 1/2 ,$$

$$\Pr\{A_1 \cap A_2\} = 1/4 ,$$

$$\Pr\{A_1 \cap A_3\} = 1/4 ,$$

$$\Pr\{A_2 \cap A_3\} = 1/4 ,$$

$$\Pr\{A_1 \cap A_2 \cap A_3\} = 0 .$$

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 (C.16) of conditional probability and the commutative law $A \cap B = B \cap A$, it follows that for two events A and B , each with nonzero probability,

$$\begin{aligned}\Pr\{A \cap B\} &= \Pr\{B\} \Pr\{A \mid B\} \\ &= \Pr\{A\} \Pr\{B \mid A\} .\end{aligned}\tag{C.18}$$

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

$$\Pr\{A \mid B\} = \frac{\Pr\{A\} \Pr\{B \mid A\}}{\Pr\{B\}} ,\tag{C.19}$$

which is known as **Bayes's theorem**. The denominator $\Pr\{B\}$ is a normalizing constant, which we can reformulate as follows. Since $B = (B \cap A) \cup (B \cap \overline{A})$, and since $B \cap A$ 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.19) produces an equivalent form of Bayes's theorem:

$$\Pr\{A \mid B\} = \frac{\Pr\{A\} \Pr\{B \mid A\}}{\Pr\{A\} \Pr\{B \mid A\} + \Pr\{\overline{A}\} \Pr\{B \mid \overline{A}\}} .\tag{C.20}$$

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

Bayes's theorem solves this problem. Let A be the event that you choose the biased coin, and let B be the event that the chosen coin comes up heads both times. We wish to determine $\Pr\{A \mid B\}$, knowing that $\Pr\{A\} = 1/2$, $\Pr\{B \mid A\} = 1$, $\Pr\{\overline{A}\} = 1/2$, and $\Pr\{B \mid \overline{A}\} = 1/4$. Thus we have

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

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. Therefore we'll assume that random variables are discrete.

For a random variable X and a real number x , we define the event $X = x$ to be $\{s \in S : X(s) = x\}$, and 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 outcomes in the sample space. Assume that the probability distribution is uniform, so that each outcome $s \in S$ 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 outcomes, namely, (1, 3), (2, 3), (3, 3), (3, 2), and (3, 1).

We can define several random variables 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.16) of conditional probability on page 1187, 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, we can define new random variables as sums, products, or other functions of the original variables.

Expected value of a random variable

The simplest, and often the 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

$$E[X] = \sum_x x \cdot \Pr\{X = x\} , \tag{C.23}$$

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

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

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

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

$$E[X + Y] = E[X] + E[Y] , \tag{C.24}$$

whenever $E[X]$ and $E[Y]$ are defined. Linearity of expectation applies to a broad range of situations, holding even when 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

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

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

$$E[aX] = aE[X] . \quad (\text{C.25})$$

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

$$E[aX + Y] = aE[X] + E[Y] . \quad (\text{C.26})$$

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

$$\begin{aligned} E[XY] &= \sum_x \sum_y xy \cdot \Pr\{X = x \text{ and } Y = y\} \\ &= \sum_x \sum_y xy \cdot \Pr\{X = x\} \Pr\{Y = y\} \quad (\text{by independence of } X \text{ and } Y) \\ &= \left(\sum_x x \cdot \Pr\{X = x\} \right) \left(\sum_y y \cdot \Pr\{Y = y\} \right) \\ &= E[X] E[Y] \quad (\text{by equation (C.23)}) . \end{aligned}$$

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

$$E[X_1 X_2 \cdots X_n] = E[X_1] E[X_2] \cdots E[X_n] . \quad (\text{C.27})$$

When a random variable X takes on values from the set of natural numbers $\mathbb{N} = \{0, 1, 2, \dots\}$, we have a nice formula for its expectation:

$$\begin{aligned} E[X] &= \sum_{i=0}^{\infty} i \cdot \Pr\{X = i\} \\ &= \sum_{i=0}^{\infty} i \cdot (\Pr\{X \geq i\} - \Pr\{X \geq i+1\}) \\ &= \sum_{i=1}^{\infty} \Pr\{X \geq i\} , \end{aligned} \quad (\text{C.28})$$

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).

A function $f(x)$ is **convex** if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (\text{C.29})$$

for all x and y and for all $0 \leq \lambda \leq 1$. **Jensen's inequality** says that when a convex function $f(x)$ is applied to a random variable X ,

$$E[f(X)] \geq f(E[X]), \quad (\text{C.30})$$

provided that the expectations exist and are finite.

Variance and standard deviation

The expected value of a random variable does not express how “spread out” the variable's values are. For example, consider 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 $E[X]$ and $E[Y]$ are $1/2$, yet the actual values taken on by Y are further 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 $E[X]$ is

$$\begin{aligned} \text{Var}[X] &= E[(X - E[X])^2] \\ &= E[X^2 - 2XE[X] + E^2[X]] \\ &= E[X^2] - 2E[XE[X]] + E^2[X] \\ &= E[X^2] - 2E^2[X] + E^2[X] \\ &= E[X^2] - E^2[X]. \end{aligned} \quad (\text{C.31})$$

To justify the equation $E[E^2[X]] = E^2[X]$, note that because $E[X]$ is a real number and not a random variable, so is $E^2[X]$. The equation $E[XE[X]] = E^2[X]$ follows from equation (C.25), with $a = E[X]$. Rewriting equation (C.31) yields an expression for the expectation of the square of a random variable:

$$E[X^2] = \text{Var}[X] + E^2[X]. \quad (\text{C.32})$$

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

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

When X and Y are independent random variables,

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

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

$$\text{Var} \left[\sum_{i=1}^n X_i \right] = \sum_{i=1}^n \text{Var} [X_i] . \quad (\text{C.33})$$

The *standard deviation* of a random variable X is the nonnegative square root of the variance of X . The standard deviation of a random variable X is sometimes denoted σ_X or simply σ when the random variable X is understood from context. With this notation, the variance of X is denoted σ^2 .

C.4 The geometric and binomial distributions

A *Bernoulli trial* is an experiment with only two possible outcomes: *success*, which occurs with probability p , and *failure*, which occurs with probability $q = 1 - p$. A coin flip serves as an example where, depending on your point of view, heads equates to success and tails to failure. 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

Consider 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 a success? Define the random variable X to be the number of trials needed to obtain a success. Then X has values in the range $\{1, 2, \dots\}$, and for $k \geq 1$,

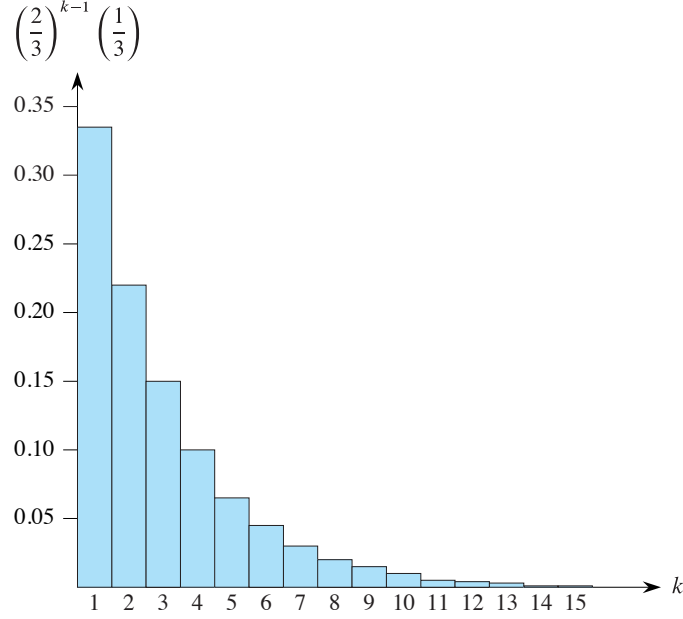


Figure C.1 A geometric distribution with probability $p = 1/3$ of success and a probability $q = 1 - p$ of failure. The expectation of the distribution is $1/p = 3$.

$$\Pr\{X = k\} = q^{k-1} p, \quad (\text{C.35})$$

since $k - 1$ failures occur before the first success. A probability distribution satisfying equation (C.35) is said to be a *geometric distribution*. Figure C.1 illustrates such a distribution.

Assuming that $q < 1$, we can calculate the expectation of a geometric distribution:

$$\begin{aligned} E[X] &= \sum_{k=1}^{\infty} k q^{k-1} p \\ &= \frac{p}{q} \sum_{k=0}^{\infty} k q^k \\ &= \frac{p}{q} \cdot \frac{q}{(1-q)^2} \quad (\text{by equation (A.11) on page 1142}) \\ &= \frac{p}{q} \cdot \frac{q}{p^2} \\ &= 1/p. \end{aligned} \quad (\text{C.36})$$

Thus, on average, it takes $1/p$ trials before a success occurs, an intuitive result. As Exercise C.4-3 asks you to show, the variance is

$$\text{Var}[X] = q/p^2. \quad (\text{C.37})$$

As an example, suppose that you repeatedly roll two dice until you 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 you'd have to roll $1/p = 9/2 = 4.5$ times on average to obtain a seven or eleven.

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 $\{0, 1, \dots, n\}$, and for $k = 0, 1, \dots, n$,

$$\Pr\{X = k\} = \binom{n}{k} p^k q^{n-k}, \quad (\text{C.38})$$

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.38) is said to be a **binomial distribution**. For convenience, we define the family of binomial distributions using the notation

$$b(k; n, p) = \binom{n}{k} p^k (1 - p)^{n-k}. \quad (\text{C.39})$$

Figure C.2 illustrates a binomial distribution. The name “binomial” comes from the right-hand side of equation (C.38) being the k th term of the expansion of $(p + q)^n$. Consequently, since $p + q = 1$, equation (C.4) on page 1181 gives

$$\sum_{k=0}^n b(k; n, p) = 1, \quad (\text{C.40})$$

as axiom 2 of the probability axioms requires.

We can compute the expectation of a random variable having a binomial distribution from equations (C.9) and (C.40). Let X be a random variable that follows the binomial distribution $b(k; n, p)$, and let $q = 1 - p$. The definition of expectation gives

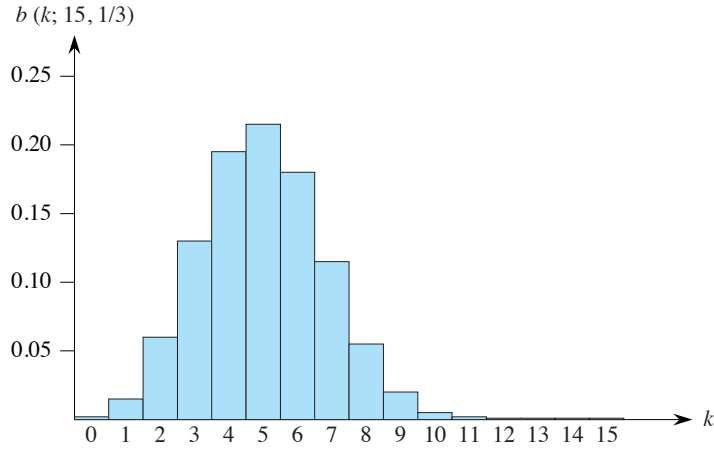


Figure C.2 The binomial distribution $b(k; 15, 1/3)$ resulting from $n = 15$ Bernoulli trials, each with probability $p = 1/3$ of success. The expectation of the distribution is $np = 5$.

$$\begin{aligned}
 E[X] &= \sum_{k=0}^n k \cdot \Pr\{X = k\} \\
 &= \sum_{k=0}^n k \cdot b(k; n, p) \\
 &= \sum_{k=1}^n k \binom{n}{k} p^k q^{n-k} \\
 &= np \sum_{k=1}^n \binom{n-1}{k-1} p^{k-1} q^{n-k} \quad (\text{by equation (C.9) on page 1183}) \\
 &= np \sum_{k=0}^{n-1} \binom{n-1}{k} p^k q^{(n-1)-k} \\
 &= np \sum_{k=0}^{n-1} b(k; n-1, p) \\
 &= np \quad (\text{by equation (C.40)}) . \tag{C.41}
 \end{aligned}$$

Linearity of expectation produces 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 $E[X_i] = p \cdot 1 + q \cdot 0 = p$, and the expected number of successes for n trials is

$$\begin{aligned}
E[X] &= E\left[\sum_{i=1}^n X_i\right] \\
&= \sum_{i=1}^n E[X_i] \quad (\text{by equation (C.24) on page 1192}) \\
&= \sum_{i=1}^n p \\
&= np.
\end{aligned} \tag{C.42}$$

We can use the same approach to calculate the variance of the distribution. By equation (C.31), $\text{Var}[X_i] = E[X_i^2] - E^2[X_i]$. Since X_i takes on only the values 0 and 1, we have $X_i^2 = X_i$, which implies $E[X_i^2] = E[X_i] = p$. Hence,

$$\text{Var}[X_i] = p - p^2 = p(1 - p) = pq. \tag{C.43}$$

To compute the variance of X , we take advantage of the independence of the n trials. By equation (C.33), we have

$$\begin{aligned}
\text{Var}[X] &= \text{Var}\left[\sum_{i=1}^n X_i\right] \\
&= \sum_{i=1}^n \text{Var}[X_i] \\
&= \sum_{i=1}^n pq \\
&= npq.
\end{aligned} \tag{C.44}$$

As Figure C.2 shows, the binomial distribution $b(k; n, p)$ increases with k until it reaches the mean np , and then it decreases. To prove that the distribution always behaves in this manner, examine the ratio of successive terms:

$$\begin{aligned}
\frac{b(k; n, p)}{b(k-1; n, p)} &= \frac{\binom{n}{k} p^k q^{n-k}}{\binom{n}{k-1} p^{k-1} q^{n-k+1}} \\
&= \frac{n! (k-1)! (n-k+1)! p}{k! (n-k)! n! q} \\
&= \frac{(n-k+1)p}{kq} \\
&= 1 + \frac{(n-k+1)p - kq}{kq} \\
&= 1 + \frac{(n-k+1)p - k(1-p)}{kq}
\end{aligned} \tag{C.45}$$

$$= 1 + \frac{(n+1)p - k}{kq}.$$

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 $(n+1)p$ is an integer, then for $k = (n+1)p$, the ratio $b(k; n, p)/b(k-1; n, p)$ equals 1, so that $b(k; n, p) = b(k-1; n, p)$. In this case, 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 $np - q < k < (n+1)p$.