

# A *Probability*

*We review briefly the elements of probability, the concept of a random variable, and example distributions.*

## A.1 Elements of Probability

A **RANDOM** experiment is one whose outcome is not predictable with certainty in advance (Ross 1987; Casella and Berger 1990). The set of all possible outcomes is known as the *sample space*  $S$ . A sample space is *discrete* if it consists of a finite (or countably infinite) set of outcomes; otherwise it is *continuous*. Any subset  $E$  of  $S$  is an *event*. Events are sets, and we can talk about their complement, intersection, union, and so forth.

One interpretation of probability is as a *frequency*. When an experiment is continually repeated under the exact same conditions, for any event  $E$ , the proportion of time that the outcome is in  $E$  approaches some constant value. This constant limiting frequency is the probability of the event, and we denote it as  $P(E)$ .

Probability sometimes is interpreted as a *degree of belief*. For example, when we speak of Turkey's probability of winning the World Soccer Cup in 2018, we do not mean a frequency of occurrence, since the championship will happen only once and it has not yet occurred (at the time of the writing of this book). What we mean in such a case is a subjective degree of belief in the occurrence of the event. Because it is subjective, different individuals may assign different probabilities to the same event.

### A.1.1 Axioms of Probability

Axioms ensure that the probabilities assigned in a random experiment can be interpreted as relative frequencies and that the assignments are consistent with our intuitive understanding of relationships among relative frequencies:

1.  $0 \leq P(E) \leq 1$ . If  $E_1$  is an event that cannot possibly occur, then  $P(E_1) = 0$ . If  $E_2$  is sure to occur,  $P(E_2) = 1$ .
2.  $S$  is the sample space containing all possible outcomes,  $P(S) = 1$ .
3. If  $E_i, i = 1, \dots, n$  are mutually exclusive (i.e., if they cannot occur at the same time, as in  $E_i \cap E_j = \emptyset, j \neq i$ , where  $\emptyset$  is the *null event* that does not contain any possible outcomes), we have

$$(A.1) \quad P\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n P(E_i)$$

For example, letting  $E^c$  denote the *complement* of  $E$ , consisting of all possible outcomes in  $S$  that are not in  $E$ , we have  $E \cap E^c = \emptyset$  and

$$\begin{aligned} P(E \cup E^c) &= P(E) + P(E^c) = 1 \\ P(E^c) &= 1 - P(E) \end{aligned}$$

If the intersection of  $E$  and  $F$  is not empty, we have

$$(A.2) \quad P(E \cup F) = P(E) + P(F) - P(E \cap F)$$

### A.1.2 Conditional Probability

$P(E|F)$  is the probability of the occurrence of event  $E$  given that  $F$  occurred and is given as

$$(A.3) \quad P(E|F) = \frac{P(E \cap F)}{P(F)}$$

Knowing that  $F$  occurred reduces the sample space to  $F$ , and the part of it where  $E$  also occurred is  $E \cap F$ . Note that equation A.3 is well-defined only if  $P(F) > 0$ . Because  $\cap$  is commutative, we have

$$P(E \cap F) = P(E|F)P(F) = P(F|E)P(E)$$

which gives us *Bayes' formula*:

$$(A.4) \quad P(F|E) = \frac{P(E|F)P(F)}{P(E)}$$

When  $F_i$  are mutually exclusive and exhaustive, namely,  $\bigcup_{i=1}^n F_i = S$

$$(A.5) \quad \begin{aligned} E &= \bigcup_{i=1}^n E \cap F_i \\ P(E) &= \sum_{i=1}^n P(E \cap F_i) = \sum_{i=1}^n P(E|F_i)P(F_i) \end{aligned}$$

Bayes' formula allows us to write

$$(A.6) \quad P(F_i|E) = \frac{P(E \cap F_i)}{P(E)} = \frac{P(E|F_i)P(F_i)}{\sum_j P(E|F_j)P(F_j)}$$

If  $E$  and  $F$  are *independent*, we have  $P(E|F) = P(E)$  and thus

$$(A.7) \quad P(E \cap F) = P(E)P(F)$$

That is, knowledge of whether  $F$  has occurred does not change the probability that  $E$  occurs.

## A.2 Random Variables

A *random variable* is a function that assigns a number to each outcome in the sample space of a random experiment.

### A.2.1 Probability Distribution and Density Functions

The *probability distribution function*  $F(\cdot)$  of a random variable  $X$  for any real number  $a$  is

$$(A.8) \quad F(a) = P\{X \leq a\}$$

and we have

$$(A.9) \quad P\{a < X \leq b\} = F(b) - F(a)$$

If  $X$  is a discrete random variable

$$(A.10) \quad F(a) = \sum_{\forall x \leq a} P(x)$$

where  $P(\cdot)$  is the *probability mass function* defined as  $P(a) = P\{X = a\}$ . If  $X$  is a *continuous* random variable,  $p(\cdot)$  is the *probability density function* such that

$$(A.11) \quad F(a) = \int_{-\infty}^a p(x) dx$$

### A.2.2 Joint Distribution and Density Functions

In certain experiments, we may be interested in the relationship between two or more random variables, and we use the *joint* probability distribution and density functions of  $X$  and  $Y$  satisfying

$$(A.12) \quad F(x, y) = P\{X \leq x, Y \leq y\}$$

Individual *marginal* distributions and densities can be computed by marginalizing, namely, summing over the free variable:

$$(A.13) \quad F_X(x) = P\{X \leq x\} = P\{X \leq x, Y \leq \infty\} = F(x, \infty)$$

In the discrete case, we write

$$(A.14) \quad P(X = x) = \sum_j P(x, y_j)$$

and in the continuous case, we have

$$(A.15) \quad p_X(x) = \int_{-\infty}^{\infty} p(x, y) dy$$

If  $X$  and  $Y$  are *independent*, we have

$$(A.16) \quad p(x, y) = p_X(x)p_Y(y)$$

These can be generalized in a straightforward manner to more than two random variables.

### A.2.3 Conditional Distributions

When  $X$  and  $Y$  are random variables

$$(A.17) \quad P_{X|Y}(x|y) = P\{X = x|Y = y\} = \frac{P\{X = x, Y = y\}}{P\{Y = y\}} = \frac{P(x, y)}{P_Y(y)}$$

### A.2.4 Bayes' Rule

When two random variables are jointly distributed with the value of one known, the probability that the other takes a given value can be computed using *Bayes' rule*:

$$(A.18) \quad P(y|x) = \frac{P(x|y)P_Y(y)}{P_X(x)} = \frac{P(x|y)P_Y(y)}{\sum_y P(x|y)P_Y(y)}$$

Or, in words

$$(A.19) \quad \text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

Note that the denominator is obtained by summing (or integrating if  $y$  is continuous) the numerator over all possible  $y$  values. The “shape” of  $p(y|x)$  depends on the numerator with denominator as a normalizing factor to guarantee that  $p(y|x)$  sum to 1. Bayes' rule allows us to modify a prior probability into a posterior probability by taking information provided by  $x$  into account.

Bayes' rule inverts dependencies, allowing us to compute  $p(y|x)$  if  $p(x|y)$  is known. Suppose that  $y$  is the “cause” of  $x$ , like  $y$  going on summer vacation and  $x$  having a suntan. Then  $p(x|y)$  is the probability that someone who is known to have gone on summer vacation has a suntan. This is the *causal* (or predictive) way. Bayes' rule allows us a *diagnostic* approach by allowing us to compute  $p(y|x)$ : namely, the probability that someone who is known to have a suntan, has gone on summer vacation. Then  $p(y)$  is the general probability of anyone's going on summer vacation and  $p(x)$  is the probability that anyone has a suntan, including both those who have gone on summer vacation and those who have not.

### A.2.5 Expectation

*Expectation*, *expected value*, or *mean* of a random variable  $X$ , denoted by  $E[X]$ , is the average value of  $X$  in a large number of experiments:

$$(A.20) \quad E[X] = \begin{cases} \sum_i x_i P(x_i) & \text{if } X \text{ is discrete} \\ \int x p(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

It is a weighted average where each value is weighted by the probability that  $X$  takes that value. It has the following properties ( $a, b \in \mathbb{R}$ ):

$$(A.21) \quad \begin{aligned} E[aX + b] &= aE[X] + b \\ E[X + Y] &= E[X] + E[Y] \end{aligned}$$

For any real-valued function  $g(\cdot)$ , the expected value is

$$(A.22) \quad E[g(X)] = \begin{cases} \sum_i g(x_i)P(x_i) & \text{if } X \text{ is discrete} \\ \int g(x)p(x)dx & \text{if } X \text{ is continuous} \end{cases}$$

A special  $g(x) = x^n$ , called the  $n$ th moment of  $X$ , is defined as

$$(A.23) \quad E[X^n] = \begin{cases} \sum_i x_i^n P(x_i) & \text{if } X \text{ is discrete} \\ \int x^n p(x)dx & \text{if } X \text{ is continuous} \end{cases}$$

*Mean* is the first moment and is denoted by  $\mu$ .

### A.2.6 Variance

*Variance* measures how much  $X$  varies around the expected value. If  $\mu \equiv E[X]$ , the variance is defined as

$$(A.24) \quad \text{Var}(X) = E[(X - \mu)^2] = E[X^2] - \mu^2$$

Variance is the second moment minus the square of the first moment. Variance, denoted by  $\sigma^2$ , satisfies the following property ( $a, b \in \mathfrak{R}$ ):

$$(A.25) \quad \text{Var}(aX + b) = a^2 \text{Var}(X)$$

$\sqrt{\text{Var}(X)}$  is called the *standard deviation* and is denoted by  $\sigma$ . Standard deviation has the same unit as  $X$  and is easier to interpret than variance.

*Covariance* indicates the relationship between two random variables. If the occurrence of  $X$  makes  $Y$  more likely to occur, then the covariance is positive; it is negative if  $X$ 's occurrence makes  $Y$  less likely to happen and is 0 if there is no dependence.

$$(A.26) \quad \text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = E[XY] - \mu_X \mu_Y$$

where  $\mu_X \equiv E[X]$  and  $\mu_Y \equiv E[Y]$ . Some other properties are

$$(A.27) \quad \begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(Y, X) \\ \text{Cov}(X, X) &= \text{Var}(X) \\ \text{Cov}(X + Z, Y) &= \text{Cov}(X, Y) + \text{Cov}(Z, Y) \\ \text{Cov}\left(\sum_i X_i, Y\right) &= \sum_i \text{Cov}(X_i, Y) \end{aligned}$$

$$(A.28) \quad \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

$$(A.29) \quad \text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i) + \sum_i \sum_{j \neq i} \text{Cov}(X_i, X_j)$$

If  $X$  and  $Y$  are independent,  $E[XY] = E[X]E[Y] = \mu_X\mu_Y$  and  $\text{Cov}(X, Y) = 0$ . Thus if  $X_i$  are independent

$$(A.30) \quad \text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i)$$

*Correlation* is a normalized, dimensionless quantity that is always between  $-1$  and  $1$ :

$$(A.31) \quad \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

### A.2.7 Weak Law of Large Numbers

Let  $\mathcal{X} = \{X^t\}_{t=1}^N$  be a set of independent and identically distributed (iid) random variables each having mean  $\mu$  and a finite variance  $\sigma^2$ . Then for any  $\epsilon > 0$ ,

$$(A.32) \quad P\left\{\left|\frac{\sum_t X^t}{N} - \mu\right| > \epsilon\right\} \rightarrow 0 \text{ as } N \rightarrow \infty$$

That is, the average of  $N$  trials converges to the mean as  $N$  increases.

## A.3 Special Random Variables

There are certain types of random variables that occur so frequently that names are given to them.

### A.3.1 Bernoulli Distribution

A trial is performed whose outcome is either a “success” or a “failure.” The random variable  $X$  is a 0/1 indicator variable and takes the value 1 for a success outcome and is 0 otherwise.  $p$  is the probability that the result of trial is a success. Then

$$(A.33) \quad P\{X = 1\} = p \text{ and } P\{X = 0\} = 1 - p$$

which can equivalently be written as

$$(A.34) \quad P\{X = i\} = p^i(1 - p)^{1-i}, i = 0, 1$$

If  $X$  is Bernoulli, its expected value and variance are

$$(A.35) \quad E[X] = p, \text{ Var}(X) = p(1 - p)$$

### A.3.2 Binomial Distribution

If  $N$  identical independent Bernoulli trials are made, the random variable  $X$  that represents the number of successes that occurs in  $N$  trials is binomial distributed. The probability that there are  $i$  successes is

$$(A.36) \quad P\{X = i\} = \binom{N}{i} p^i (1-p)^{N-i}, i = 0 \dots N$$

If  $X$  is binomial, its expected value and variance are

$$(A.37) \quad E[X] = Np, \text{Var}(X) = Np(1-p)$$

### A.3.3 Multinomial Distribution

Consider a generalization of Bernoulli where instead of two states, the outcome of a random event is one of  $K$  mutually exclusive and exhaustive states, each of which has a probability of occurring  $p_i$  where  $\sum_{i=1}^K p_i = 1$ . Suppose that  $N$  such trials are made where outcome  $i$  occurred  $N_i$  times with  $\sum_{i=1}^K N_i = N$ . Then the joint distribution of  $N_1, N_2, \dots, N_K$  is multinomial:

$$(A.38) \quad P(N_1, N_2, \dots, N_K) = N! \prod_{i=1}^K \frac{p_i^{N_i}}{N_i!}$$

A special case is when  $N = 1$ ; only one trial is made. Then  $N_i$  are 0/1 indicator variables of which only one of them is 1 and all others are 0. Then equation A.38 reduces to

$$(A.39) \quad P(N_1, N_2, \dots, N_K) = \prod_{i=1}^K p_i^{N_i}$$

### A.3.4 Uniform Distribution

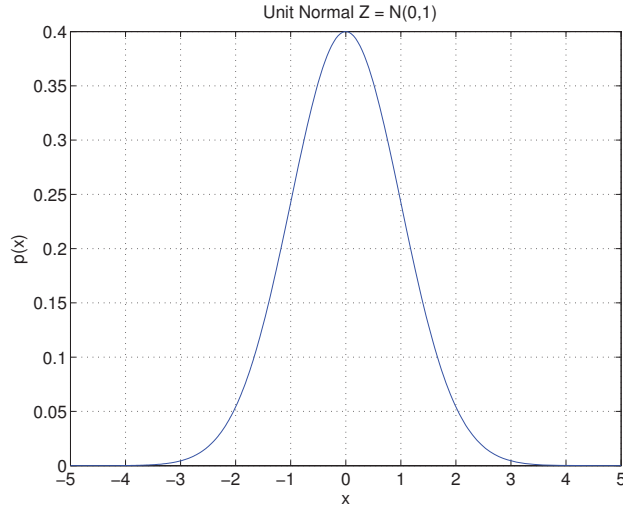
$X$  is uniformly distributed over the interval  $[a, b]$  if its density function is given by

$$(A.40) \quad p(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

If  $X$  is uniform, its expected value and variance are

$$(A.41) \quad E[X] = \frac{a+b}{2}, \text{Var}(X) = \frac{(b-a)^2}{12}$$





**Figure A.1** Probability density function of  $Z$ , the unit normal distribution.

### A.3.5 Normal (Gaussian) Distribution

$X$  is normal or Gaussian distributed with mean  $\mu$  and variance  $\sigma^2$ , denoted as  $\mathcal{N}(\mu, \sigma^2)$ , if its density function is

$$(A.42) \quad p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad -\infty < x < \infty$$

Many random phenomena obey the bell-shaped normal distribution, at least approximately, and many observations from nature can be seen as a continuous, slightly different versions of a typical value—that is probably why it is called the *normal* distribution. In such a case,  $\mu$  represents the typical value and  $\sigma$  defines how much instances vary around the prototypical value.

68.27 percent lie in  $(\mu - \sigma, \mu + \sigma)$ , 95.45 percent in  $(\mu - 2\sigma, \mu + 2\sigma)$ , and 99.73 percent in  $(\mu - 3\sigma, \mu + 3\sigma)$ . Thus  $P\{|x - \mu| < 3\sigma\} \approx 0.99$ . For practical purposes,  $p(x) \approx 0$  if  $x < \mu - 3\sigma$  or  $x > \mu + 3\sigma$ .  $Z$  is unit normal, namely,  $\mathcal{N}(0, 1)$  (see figure A.1), and its density is written as

$$(A.43) \quad p_Z(x) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{x^2}{2}\right]$$

If  $X \sim \mathcal{N}(\mu, \sigma^2)$  and  $Y = aX + b$ , then  $Y \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$ . The sum of independent normal variables is also normal with  $\mu = \sum_i \mu_i$  and  $\sigma^2 = \sum_i \sigma_i^2$ . If  $X$  is  $\mathcal{N}(\mu, \sigma^2)$ , then

$$(A.44) \quad \frac{X - \mu}{\sigma} \sim \mathcal{Z}$$

This is called z-normalization.

CENTRAL LIMIT  
THEOREM Let  $X_1, X_2, \dots, X_N$  be a set of iid random variables all having mean  $\mu$  and variance  $\sigma^2$ . Then the *central limit theorem* states that for large  $N$ , the distribution of

$$(A.45) \quad X_1 + X_2 + \dots + X_N$$

is approximately  $\mathcal{N}(N\mu, N\sigma^2)$ . For example, if  $X$  is binomial with parameters  $(N, p)$ ,  $X$  can be written as the sum of  $N$  Bernoulli trials and  $(X - Np)/\sqrt{Np(1-p)}$  is approximately unit normal.

Central limit theorem is also used to generate normally distributed random variables on computers. Programming languages have subroutines that return uniformly distributed (pseudo-)random numbers in the range  $[0, 1]$ . When  $U_i$  are such random variables,  $\sum_{i=1}^{12} U_i - 6$  is approximately  $\mathcal{Z}$ .

Let us say  $X^t \sim \mathcal{N}(\mu, \sigma^2)$ . The estimated sample mean

$$(A.46) \quad m = \frac{\sum_{t=1}^N X^t}{N}$$

is also normal with mean  $\mu$  and variance  $\sigma^2/N$ .

### A.3.6 Chi-Square Distribution

If  $Z_i$  are independent unit normal random variables, then

$$(A.47) \quad X = Z_1^2 + Z_2^2 + \dots + Z_n^2$$

is chi-square with  $n$  degrees of freedom, namely,  $X \sim \mathcal{X}_n^2$ , with

$$(A.48) \quad E[X] = n, \text{ Var}(X) = 2n$$

When  $X^t \sim \mathcal{N}(\mu, \sigma^2)$ , the estimated sample variance is

$$(A.49) \quad S^2 = \frac{\sum_t (X^t - m)^2}{N - 1}$$

and we have

$$(A.50) \quad (N - 1) \frac{S^2}{\sigma^2} \sim \mathcal{X}_{N-1}^2$$

It is also known that  $m$  and  $S^2$  are independent.

**A.3.7  $t$  Distribution**

If  $Z \sim \mathcal{Z}$  and  $X \sim \chi_n^2$  are independent, then

$$(A.51) \quad T_n = \frac{Z}{\sqrt{X/n}}$$

is  $t$ -distributed with  $n$  degrees of freedom with

$$(A.52) \quad E[T_n] = 0, n > 1, \text{Var}(T_n) = \frac{n}{n-2}, n > 2$$

Like the unit normal density,  $t$  is symmetric around 0. As  $n$  becomes larger,  $t$  density becomes more and more like the unit normal, the difference being that  $t$  has thicker tails, indicating greater variability than does normal.

**A.3.8  $F$  Distribution**

If  $X_1 \sim \chi_n^2$  and  $X_2 \sim \chi_m^2$  are independent chi-square random variables with  $n$  and  $m$  degrees of freedom, respectively,

$$(A.53) \quad F_{n,m} = \frac{X_1/n}{X_2/m}$$

is  $F$ -distributed with  $n$  and  $m$  degrees of freedom with

$$(A.54) \quad E[F_{n,m}] = \frac{m}{m-2}, m > 2, \text{Var}(F_{n,m}) = \frac{m^2(2m+2n-4)}{n(m-2)^2(m-4)}, m > 4$$

**A.4 References**

- Casella, G., and R. L. Berger. 1990. *Statistical Inference*. Belmont, CA: Duxbury.
- Ross, S. M. 1987. *Introduction to Probability and Statistics for Engineers and Scientists*. New York: Wiley.