

We can represent this system as the matrix equation  $\mathbf{A} \mathbf{x} = \mathbf{b}$ , where

$$\mathbf{A} = \begin{pmatrix} 2 & 1 & -1 \\ -3 & -1 & 2 \\ -2 & 1 & 2 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x \\ y \\ z \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 8 \\ -11 \\ -3 \end{pmatrix}.$$

To solve  $\mathbf{A} \mathbf{x} = \mathbf{b}$  we multiply both sides by  $\mathbf{A}^{-1}$ , yielding  $\mathbf{A}^{-1} \mathbf{A} \mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$ , which simplifies to  $\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$ . After inverting  $\mathbf{A}$  and multiplying by  $\mathbf{b}$ , we get the answer

$$\mathbf{x} = \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 2 \\ 3 \\ -1 \end{pmatrix}.$$

## A.3 PROBABILITY DISTRIBUTIONS

---

A probability is a measure over a set of events that satisfies three axioms:

1. The measure of each event is between 0 and 1. We write this as  $0 \leq P(X = x_i) \leq 1$ , where  $X$  is a random variable representing an event and  $x_i$  are the possible values of  $X$ . In general, random variables are denoted by uppercase letters and their values by lowercase letters.
2. The measure of the whole set is 1; that is,  $\sum_{i=1}^n P(X = x_i) = 1$ .
3. The probability of a union of disjoint events is the sum of the probabilities of the individual events; that is,  $P(X = x_1 \vee X = x_2) = P(X = x_1) + P(X = x_2)$ , where  $x_1$  and  $x_2$  are disjoint.

A **probabilistic model** consists of a sample space of mutually exclusive possible outcomes, together with a probability measure for each outcome. For example, in a model of the weather tomorrow, the outcomes might be *sunny*, *cloudy*, *rainy*, and *snowy*. A subset of these outcomes constitutes an event. For example, the event of precipitation is the subset consisting of  $\{\text{rainy, snowy}\}$ .

We use  $\mathbf{P}(X)$  to denote the vector of values  $\langle P(X = x_1), \dots, P(X = x_n) \rangle$ . We also use  $P(x_i)$  as an abbreviation for  $P(X = x_i)$  and  $\sum_x P(x)$  for  $\sum_{i=1}^n P(X = x_i)$ .

The conditional probability  $P(B|A)$  is defined as  $P(B \cap A)/P(A)$ .  $A$  and  $B$  are conditionally independent if  $P(B|A) = P(B)$  (or equivalently,  $P(A|B) = P(A)$ ). For continuous variables, there are an infinite number of values, and unless there are point spikes, the probability of any one value is 0. Therefore, we define a **probability density function**, which we also denote as  $P(\cdot)$ , but which has a slightly different meaning from the discrete probability function. The density function  $P(x)$  for a random variable  $X$ , which might be thought of as  $P(X = x)$ , is intuitively defined as the ratio of the probability that  $X$  falls into an interval around  $x$ , divided by the width of the interval, as the interval width goes to zero:

$$P(x) = \lim_{dx \rightarrow 0} P(x \leq X \leq x + dx)/dx.$$

The density function must be nonnegative for all  $x$  and must have

$$\int_{-\infty}^{\infty} P(x) dx = 1.$$

CUMULATIVE  
PROBABILITY  
DENSITY FUNCTION

We can also define a **cumulative probability density function**  $F_X(x)$ , which is the probability of a random variable being less than  $x$ :

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x P(u) du.$$

Note that the probability density function has units, whereas the discrete probability function is unitless. For example, if values of  $X$  are measured in seconds, then the density is measured in Hz (i.e., 1/sec). If values of  $\mathbf{X}$  are points in three-dimensional space measured in meters, then density is measured in  $1/m^3$ .

GAUSSIAN  
DISTRIBUTION

One of the most important probability distributions is the **Gaussian distribution**, also known as the **normal distribution**. A Gaussian distribution with mean  $\mu$  and standard deviation  $\sigma$  (and therefore variance  $\sigma^2$ ) is defined as

$$P(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/(2\sigma^2)},$$

STANDARD NORMAL  
DISTRIBUTION  
MULTIVARIATE  
GAUSSIAN

where  $x$  is a continuous variable ranging from  $-\infty$  to  $+\infty$ . With mean  $\mu=0$  and variance  $\sigma^2=1$ , we get the special case of the **standard normal distribution**. For a distribution over a vector  $\mathbf{x}$  in  $n$  dimensions, there is the **multivariate Gaussian distribution**:

$$P(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}((\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}-\boldsymbol{\mu}))},$$

CUMULATIVE  
DISTRIBUTION

where  $\boldsymbol{\mu}$  is the mean vector and  $\Sigma$  is the **covariance matrix** (see below).

In one dimension, we can define the **cumulative distribution** function  $F(x)$  as the probability that a random variable will be less than  $x$ . For the normal distribution, this is

$$F(x) = \int_{-\infty}^x P(z) dz = \frac{1}{2} (1 + \text{erf}(\frac{z-\mu}{\sigma\sqrt{2}})),$$

CENTRAL LIMIT  
THEOREM

where  $\text{erf}(x)$  is the so-called **error function**, which has no closed-form representation.

The **central limit theorem** states that the distribution formed by sampling  $n$  independent random variables and taking their mean tends to a normal distribution as  $n$  tends to infinity. This holds for almost any collection of random variables, even if they are not strictly independent, unless the variance of any finite subset of variables dominates the others.

EXPECTATION

The **expectation** of a random variable,  $E(X)$ , is the mean or average value, weighted by the probability of each value. For a discrete variable it is:

$$E(X) = \sum_i x_i P(X = x_i).$$

For a continuous variable, replace the summation with an integral over the probability density function,  $P(x)$ :

$$E(X) = \int_{-\infty}^{\infty} x P(x) dx,$$

ROOT MEAN SQUARE

The **root mean square**, RMS, of a set of values (often samples of a random variable) is the square root of the mean of the squares of the values,

$$RMS(x_1, \dots, x_n) = \sqrt{\frac{x_1^2 + \dots + x_n^2}{n}}.$$

COVARIANCE

The **covariance** of two random variables is the expectation of the product of their differences from their means:

$$\text{cov}(X, Y) = E((X - \mu_X)(Y - \mu_Y)).$$

COVARIANCE MATRIX

The **covariance matrix**, often denoted  $\Sigma$ , is a matrix of covariances between elements of a vector of random variables. Given  $\mathbf{X} = \langle X_1, \dots, X_n \rangle^\top$ , the entries of the covariance matrix are as follows:

$$\Sigma_{i,j} = \text{cov}(X_i, X_j) = E((X_i - \mu_i)(X_j - \mu_j)).$$

A few more miscellaneous points: we use  $\log(x)$  for the natural logarithm,  $\log_e(x)$ . We use  $\text{argmax}_x f(x)$  for the value of  $x$  for which  $f(x)$  is maximal.

## BIBLIOGRAPHICAL AND HISTORICAL NOTES

The  $O()$  notation so widely used in computer science today was first introduced in the context of number theory by the German mathematician P. G. H. Bachmann (1894). The concept of NP-completeness was invented by Cook (1971), and the modern method for establishing a reduction from one problem to another is due to Karp (1972). Cook and Karp have both won the Turing award, the highest honor in computer science, for their work.

Classic works on the analysis and design of algorithms include those by Knuth (1973) and Aho, Hopcroft, and Ullman (1974); more recent contributions are by Tarjan (1983) and Cormen, Leiserson, and Rivest (1990). These books place an emphasis on designing and analyzing algorithms to solve tractable problems. For the theory of NP-completeness and other forms of intractability, see Garey and Johnson (1979) or Papadimitriou (1994). Good texts on probability include Chung (1979), Ross (1988), and Bertsekas and Tsitsiklis (2008).