# Appendix B

# Probability

## B.1 Foundations

The set $S$ of all possible outcomes of an experiment is called the **sample space** for the experiment. Take the simple example of tossing a coin. There are two outcomes, heads and tails, so we can write $S = \{H, T\}$. If two coins are tossed in sequence, we can write the four outcomes as $S = \{HH, HT, TH, TT\}$.

An **event** $A$ is any collection of possible outcomes of an experiment. An event is a subset of $S$, including $S$ itself and the null set $\emptyset$. Continuing the two coin example, one event is $A = \{HH, HT\}$, the event that the first coin is heads. We say that $A$ and $B$ are **disjoint** or **mutually exclusive** if $A \cap B = \emptyset$. For example, the sets $\{HH, HT\}$ and $\{TH\}$ are disjoint. Furthermore, if the sets $A_1, A_2, ...$ are pairwise disjoint and $\cup_{i=1}^{\infty} A_i = S$, then the collection $A_1, A_2, ...$ is called a **partition** of $S$.

The following are elementary set operations:

**Union:** $A \cup B = \{x : x \in A \text{ or } x \in B\}$.

**Intersection:** $A \cap B = \{x : x \in A \text{ and } x \in B\}$.

**Complement:** $A^c = \{x : x \notin A\}$.

The following are useful properties of set operations.

**Communtatitivity:** $A \cup B = B \cup A$; $\quad A \cap B = B \cap A$.

**Associativity:** $A \cup (B \cup C) = (A \cup B) \cup C$; $\quad A \cap (B \cap C) = (A \cap B) \cap C$.

**Distributive Laws:** $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$; $\quad A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$.

**DeMorgan's Laws:** $(A \cup B)^c = A^c \cap B^c$; $\quad (A \cap B)^c = A^c \cup B^c$.

A **probability function** assigns probabilities (numbers between 0 and 1) to events $A$ in $S$. This is straightforward when $S$ is countable; when $S$ is uncountable we must be somewhat more careful. A set $\mathcal{B}$ is called a **sigma algebra** (or Borel field) if $\emptyset \in \mathcal{B}$, $A \in \mathcal{B}$ implies $A^c \in \mathcal{B}$, and $A_1, A_2, ... \in \mathcal{B}$ implies $\cup_{i=1}^{\infty} A_i \in \mathcal{B}$. A simple example is $\{\emptyset, S\}$ which is known as the trivial sigma algebra. For any sample space $S$, let $\mathcal{B}$ be the smallest sigma algebra which contains all of the open sets in $S$. When $S$ is countable, $\mathcal{B}$ is simply the collection of all subsets of $S$, including $\emptyset$ and $S$. When $S$ is the real line, then $\mathcal{B}$ is the collection of all open and closed intervals. We call $\mathcal{B}$ the sigma algebra associated with $S$. We only define probabilities for events contained in $\mathcal{B}$.

We now can give the axiomatic definition of probability. Given $S$ and $\mathcal{B}$, a probability function $\mathbb{P}$ satisfies $\mathbb{P}(S) = 1$, $\mathbb{P}(A) \geq 0$ for all $A \in \mathcal{B}$, and if $A_1, A_2, ... \in \mathcal{B}$ are pairwise disjoint, then $\mathbb{P}\left(\cup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mathbb{P}(A_i)$.

Some important properties of the probability function include the following

- $\mathbb{P}(\emptyset) = 0$

- $\mathbb{P}(A) \leq 1$

- $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$

- $\mathbb{P}(B \cap A^c) = \mathbb{P}(B) - \mathbb{P}(A \cap B)$

- $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$

- If $A \subset B$ then $\mathbb{P}(A) \leq \mathbb{P}(B)$

- Bonferroni's Inequality: $\mathbb{P}(A \cap B) \geq \mathbb{P}(A) + \mathbb{P}(B) - 1$

- Boole's Inequality: $\mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B)$

For some elementary probability models, it is useful to have simple rules to count the number of objects in a set. These counting rules are facilitated by using the binomial coefficients which are defined for nonnegative integers $n$ and $r$, $n \geq r$, as

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

When counting the number of objects in a set, there are two important distinctions. Counting may be **with replacement** or **without replacement**. Counting may be **ordered** or **unordered**. For example, consider a lottery where you pick six numbers from the set 1, 2, ..., 49. This selection is without replacement if you are not allowed to select the same number twice, and is with replacement if this is allowed. Counting is ordered or not depending on whether the sequential order of the numbers is relevant to winning the lottery. Depending on these two distinctions, we have four expressions for the number of objects (possible arrangements) of size $r$ from $n$ objects.

|  | Without Replacement | With Replacement |
|---|---|---|
| Ordered | $\frac{n!}{(n-r)!}$ | $n^r$ |
| Unordered | $\binom{n}{r}$ | $\binom{n+r-1}{r}$ |

In the lottery example, if counting is unordered and without replacement, the number of potential combinations is $\binom{49}{6} = 13{,}983{,}816$.

If $\mathbb{P}(B) > 0$ the **conditional probability** of the event $A$ given the event $B$ is

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

For any $B$, the conditional probability function is a valid probability function where $S$ has been replaced by $B$. Rearranging the definition, we can write

$$\mathbb{P}(A \cap B) = \mathbb{P}(A \mid B)\,\mathbb{P}(B)$$

which is often quite useful. We can say that the occurrence of $B$ has no information about the likelihood of event $A$ when $\mathbb{P}(A \mid B) = \mathbb{P}(A)$, in which case we find

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\,\mathbb{P}(B) \tag{B.1}$$

We say that the events $A$ and $B$ are **statistically independent** when (B.1) holds. Furthermore, we say that the collection of events $A_1, ..., A_k$ are **mutually independent** when for any subset $\{A_i : i \in I\}$,

$$\mathbb{P}\left(\bigcap_{i \in I} A_i\right) = \prod_{i \in I} \mathbb{P}(A_i).$$

**Theorem 1** (Bayes' Rule). *For any set $B$ and any partition $A_1, A_2, ...$ of the sample space, then for each $i = 1, 2, ...$*

$$\mathbb{P}(A_i \mid B) = \frac{\mathbb{P}(B \mid A_i)\,\mathbb{P}(A_i)}{\sum_{j=1}^{\infty} \mathbb{P}(B \mid A_j)\,\mathbb{P}(A_j)}$$

## B.2 Random Variables

A **random variable** $X$ is a function from a sample space $S$ into the real line. This induces a new sample space – the real line – and a new probability function on the real line. Typically, we denote random variables by uppercase letters such as $X$, and use lower case letters such as $x$ for potential values and realized values. (This is in contrast to the notation adopted for most of the textbook.) For a random variable $X$ we define its **cumulative distribution function** (CDF) as

$$F(x) = \mathbb{P}\left(X \leq x\right). \tag{B.2}$$

Sometimes we write this as $F_X(x)$ to denote that it is the CDF of $X$. A function $F(x)$ is a CDF if and only if the following three properties hold:

1. $\lim_{x \to -\infty} F(x) = 0$ and $\lim_{x \to \infty} F(x) = 1$

2. $F(x)$ is nondecreasing in $x$

3. $F(x)$ is right-continuous

We say that the random variable $X$ is **discrete** if $F(x)$ is a step function. In the latter case, the range of $X$ consists of a countable set of real numbers $\tau_1, ..., \tau_r$. The probability function for $X$ takes the form

$$\mathbb{P}\left(X = \tau_j\right) = \pi_j, \qquad j = 1, ..., r \tag{B.3}$$

where $0 \leq \pi_j \leq 1$ and $\sum_{j=1}^{r} \pi_j = 1$.

We say that the random variable $X$ is **continuous** if $F(x)$ is continuous in $x$. In this case $\mathbb{P}(X = \tau) = 0$ for all $\tau \in R$ so the representation (B.3) is unavailable. Instead, we represent the relative probabilities by the **probability density function** (PDF)

$$f(x) = \frac{d}{dx} F(x)$$

so that

$$F(x) = \int_{-\infty}^{x} f(u) du$$

and

$$\mathbb{P}\left(a \leq X \leq b\right) = \int_{a}^{b} f(u) du.$$

These expressions only make sense if $F(x)$ is differentiable. While there are examples of continuous random variables which do not possess a PDF, these cases are unusual and are typically ignored.

A function $f(x)$ is a PDF if and only if $f(x) \geq 0$ for all $x \in R$ and $\int_{-\infty}^{\infty} f(x) dx$.

## B.3 Expectation

For any measurable real function $g$, we define the **mean** or **expectation** $\mathbb{E}g(X)$ as follows. If $X$ is discrete,

$$\mathbb{E}g(X) = \sum_{j=1}^{r} g(\tau_j) \pi_j,$$

and if $X$ is continuous

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

The latter is well defined and finite if

$$\int_{-\infty}^{\infty} |g(x)| f(x) dx < \infty. \tag{B.4}$$

If (B.4) does not hold, evaluate

$$I_1 = \int_{g(x)>0} g(x)f(x)dx$$

$$I_2 = -\int_{g(x)<0} g(x)f(x)dx$$

If $I_1 = \infty$ and $I_2 < \infty$ then we define $\mathbb{E}g(X) = \infty$. If $I_1 < \infty$ and $I_2 = \infty$ then we define $\mathbb{E}g(X) = -\infty$. If both $I_1 = \infty$ and $I_2 = \infty$ then $\mathbb{E}g(X)$ is undefined.

Since $\mathbb{E}(a + bX) = a + b\mathbb{E}X$, we say that expectation is a linear operator.

For $m > 0$, we define the $m'th$ **moment** of $X$ as $\mathbb{E}X^m$ and the $m'th$ **central moment** as $\mathbb{E}(X - \mathbb{E}X)^m$.

Two special moments are the **mean** $\mu = \mathbb{E}X$ and **variance** $\sigma^2 = \mathbb{E}(X - \mu)^2 = \mathbb{E}X^2 - \mu^2$. We call $\sigma = \sqrt{\sigma^2}$ the **standard deviation** of $X$. We can also write $\sigma^2 = \text{var}(X)$. For example, this allows the convenient expression $\text{var}(a + bX) = b^2\,\text{var}(X)$.

The **moment generating function** (MGF) of $X$ is

$$M(\lambda) = \mathbb{E}\exp(\lambda X).$$

The MGF does not necessarily exist. However, when it does and $\mathbb{E}|X|^m < \infty$ then

$$\left.\frac{d^m}{d\lambda^m}M(\lambda)\right|_{\lambda=0} = \mathbb{E}(X^m)$$

which is why it is called the moment generating function.

More generally, the **characteristic function** (CF) of $X$ is

$$C(\lambda) = \mathbb{E}\exp(i\lambda X)$$

where $i = \sqrt{-1}$ is the imaginary unit. The CF always exists, and when $\mathbb{E}|X|^m < \infty$

$$\left.\frac{d^m}{d\lambda^m}C(\lambda)\right|_{\lambda=0} = i^m\mathbb{E}(X^m).$$

The $L^p$ **norm**, $p \geq 1$, of the random variable $X$ is

$$\|X\|_p = (\mathbb{E}|X|^p)^{1/p}.$$

## B.4 Gamma Function

The gamma function is defined for $\alpha > 0$ as

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1}\exp(-x).$$

It satisfies the property

$$\Gamma(1 + \alpha) = \Gamma(\alpha)\alpha$$

so for positive integers $n$,

$$\Gamma(n) = (n - 1)!$$

Special values include

$$\Gamma(1) = 1$$

and

$$\Gamma\left(\frac{1}{2}\right) = \pi^{1/2}.$$

Sterling's formula is an expansion for the its logarithm

$$\log\Gamma(\alpha) = \frac{1}{2}\log(2\pi) + \left(\alpha - \frac{1}{2}\right)\log\alpha - z + \frac{1}{12\alpha} - \frac{1}{360\alpha^3} + \frac{1}{1260\alpha^5} + \cdots$$

# B.5   Common Distributions

For reference, we now list some important discrete distribution function.

**Bernoulli**

$$
\begin{aligned}
\mathbb{P}\left(X = x\right) &= p^x(1-p)^{1-x}, \qquad x = 0, 1; \qquad 0 \le p \le 1 \\
\mathbb{E}X &= p \\
\mathrm{var}(X) &= p(1-p)
\end{aligned}
$$

**Binomial**

$$
\begin{aligned}
\mathbb{P}\left(X = x\right) &= \binom{n}{x} p^x(1-p)^{n-x}, \qquad x = 0, 1, ..., n; \qquad 0 \le p \le 1 \\
\mathbb{E}X &= np \\
\mathrm{var}(X) &= np(1-p)
\end{aligned}
$$

**Geometric**

$$
\begin{aligned}
\mathbb{P}\left(X = x\right) &= p(1-p)^{x-1}, \qquad x = 1, 2, ...; \qquad 0 \le p \le 1 \\
\mathbb{E}X &= \frac{1}{p} \\
\mathrm{var}(X) &= \frac{1-p}{p^2}
\end{aligned}
$$

**Multinomial**

$$
\begin{aligned}
\mathbb{P}\left(X_1 = x_1, X_2 = x_2, ..., X_m = x_m\right) &= \frac{n!}{x_1! x_2! \cdots x_m!} p_1^{x_1} p_2^{x_2} \cdots p_m^{x_m}, \\
x_1 + \cdots + x_m &= n; \\
p_1 + \cdots + p_m &= 1 \\
\mathbb{E}X_i &= p_i \\
\mathrm{var}(X_i) &= np_i(1-p_i) \\
\mathrm{cov}\left(X_i, X_j\right) &= -np_i p_j
\end{aligned}
$$

**Negative Binomial**

$$
\begin{aligned}
\mathbb{P}\left(X = x\right) &= \frac{\Gamma\left(r+x\right)}{x! \Gamma\left(r\right)} p^r (1-p)^{x-1}, \qquad x = 0, 1, 2, ...; \qquad 0 \le p \le 1 \\
\mathbb{E}X &= \frac{r\left(1-p\right)}{p} \\
\mathrm{var}(X) &= \frac{r\left(1-p\right)}{p^2}
\end{aligned}
$$

**Poisson**

$$
\begin{aligned}
\mathbb{P}\left(X = x\right) &= \frac{\exp\left(-\lambda\right)\lambda^x}{x!}, \qquad x = 0, 1, 2, ..., \qquad \lambda > 0 \\
\mathbb{E}X &= \lambda \\
\mathrm{var}(X) &= \lambda
\end{aligned}
$$

We now list some important continuous distributions.

**Beta**

$$f(x) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}x^{\alpha-1}(1-x)^{\beta-1}, \qquad 0 \le x \le 1; \qquad \alpha > 0,\ \beta > 0$$

$$\mu = \frac{\alpha}{\alpha+\beta}$$

$$\mathrm{var}(X) = \frac{\alpha\beta}{(\alpha+\beta+1)(\alpha+\beta)^2}$$

**Cauchy**

$$f(x) = \frac{1}{\pi(1+x^2)}, \qquad -\infty < x < \infty$$

$$\mathbb{E}X = \infty$$

$$\mathrm{var}(X) = \infty$$

**Exponential**

$$f(x) = \frac{1}{\theta}\exp\left(\frac{x}{\theta}\right), \qquad 0 \le x < \infty; \qquad \theta > 0$$

$$\mathbb{E}X = \theta$$

$$\mathrm{var}(X) = \theta^2$$

**Logistic**

$$f(x) = \frac{\exp(-x)}{(1+\exp(-x))^2}, \qquad -\infty < x < \infty;$$

$$\mathbb{E}X = 0$$

$$\mathrm{var}(X) = \frac{\pi^2}{3}$$

**Lognormal**

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma x}\exp\left(-\frac{(\log x - \mu)^2}{2\sigma^2}\right), \qquad 0 \le x < \infty; \qquad \sigma > 0$$

$$\mathbb{E}X = \exp(\mu + \sigma^2/2)$$

$$\mathrm{var}(X) = \exp(2\mu + 2\sigma^2) - \exp(2\mu + \sigma^2)$$

**Pareto**

$$f(x) = \frac{\beta\alpha^\beta}{x^{\beta+1}}, \qquad \alpha \le x < \infty, \qquad \alpha > 0, \qquad \beta > 0$$

$$\mathbb{E}X = \frac{\beta\alpha}{\beta-1}, \qquad \beta > 1$$

$$\mathrm{var}(X) = \frac{\beta\alpha^2}{(\beta-1)^2(\beta-2)}, \qquad \beta > 2$$

**Uniform**

$$f(x) = \frac{1}{b-a}, \qquad a \le x \le b$$

$$\mathbb{E}X = \frac{a+b}{2}$$

$$\mathrm{var}(X) = \frac{(b-a)^2}{12}$$

**Weibull**

$$f(x) = \frac{\gamma}{\beta} x^{\gamma-1} \exp\left(-\frac{x^\gamma}{\beta}\right), \qquad 0 \le x < \infty; \qquad \gamma > 0,\ \beta > 0$$

$$\mathbb{E}X = \beta^{1/\gamma} \Gamma\left(1 + \frac{1}{\gamma}\right)$$

$$\mathrm{var}(X) = \beta^{2/\gamma}\left(\Gamma\left(1 + \frac{2}{\gamma}\right) - \Gamma^2\left(1 + \frac{1}{\gamma}\right)\right)$$

**Gamma**

$$f(x) = \frac{1}{\Gamma(\alpha)\theta^\alpha} x^{\alpha-1} \exp\left(-\frac{x}{\theta}\right), \qquad 0 \le x < \infty; \qquad \alpha > 0,\ \theta > 0$$

$$\mathbb{E}X = \alpha\theta$$

$$\mathrm{var}(X) = \alpha\theta^2$$

**Chi-Square**

$$f(x) = \frac{1}{\Gamma(r/2)2^{r/2}} x^{r/2-1} \exp\left(-\frac{x}{2}\right), \qquad 0 \le x < \infty; \qquad r > 0$$

$$\mathbb{E}X = r$$

$$\mathrm{var}(X) = 2r$$

**Normal**

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \qquad -\infty < x < \infty; \qquad -\infty < \mu < \infty,\ \sigma^2 > 0$$

$$\mathbb{E}X = \mu$$

$$\mathrm{var}(X) = \sigma^2$$

**Student t**

$$f(x) = \frac{\Gamma\left(\frac{r+1}{2}\right)}{\sqrt{r\pi}\Gamma\left(\frac{r}{2}\right)} \left(1 + \frac{x^2}{r}\right)^{-\left(\frac{r+1}{2}\right)}, \qquad -\infty < x < \infty; \qquad r > 0$$

$$\mathbb{E}X = 0 \text{ if } r > 1$$

$$\mathrm{var}(X) = \frac{r}{r-2} \text{ if } r > 2$$

## B.6 Multivariate Random Variables

A pair of bivariate random variables $(X, Y)$ is a function from the sample space into $\mathbb{R}^2$. The joint CDF of $(X, Y)$ is

$$F(x, y) = \mathbb{P}\left(X \le x, Y \le y\right).$$

If $F$ is continuous, the joint probability density function is

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

For a Borel measurable set $A \in R^2$,

$$\mathbb{P}\left((X < Y) \in A\right) = \int\int_A f(x, y)dxdy$$

For any measurable function $g(x,y)$,

$$\mathbb{E}g(X,Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y)f(x,y)dxdy.$$

The **marginal distribution** of $X$ is

$$
\begin{aligned}
F_X(x) &= \mathbb{P}(X \leq x) \\
&= \lim_{y \to \infty} F(x,y) \\
&= \int_{-\infty}^{x} \int_{-\infty}^{\infty} f(x,y)dydx
\end{aligned}
$$

so the **marginal density** of $X$ is

$$f_X(x) = \frac{d}{dx}F_X(x) = \int_{-\infty}^{\infty} f(x,y)dy.$$

Similarly, the marginal density of $Y$ is

$$f_Y(y) = \int_{-\infty}^{\infty} f(x,y)dx.$$

The random variables $X$ and $Y$ are defined to be **independent** if $f(x,y) = f_X(x)f_Y(y)$. Furthermore, $X$ and $Y$ are independent if and only if there exist functions $g(x)$ and $h(y)$ such that $f(x,y) = g(x)h(y)$.

If $X$ and $Y$ are independent, then

$$
\begin{aligned}
\mathbb{E}\left(g(X)h(Y)\right) &= \int \int g(x)h(y)f(y,x)dydx \\
&= \int \int g(x)h(y)f_Y(y)f_X(x)dydx \\
&= \int g(x)f_X(x)dx \int h(y)f_Y(y)dy \\
&= \mathbb{E}g(X)\,\mathbb{E}h(Y).
\end{aligned}
\tag{B.5}
$$

if the expectations exist. For example, if $X$ and $Y$ are independent then

$$\mathbb{E}(XY) = \mathbb{E}X\mathbb{E}Y.$$

Another implication of (B.5) is that if $X$ and $Y$ are independent and $Z = X + Y$, then

$$
\begin{aligned}
M_Z(\lambda) &= \mathbb{E}\exp\left(\lambda\left(X+Y\right)\right) \\
&= \mathbb{E}\left(\exp\left(\lambda X\right)\exp\left(\lambda Y\right)\right) \\
&= \mathbb{E}\exp\left(\lambda'X\right)\mathbb{E}\exp\left(\lambda'Y\right) \\
&= M_X(\lambda)M_Y(\lambda).
\end{aligned}
\tag{B.6}
$$

The covariance between $X$ and $Y$ is

$$\text{cov}(X,Y) = \sigma_{XY} = \mathbb{E}\left(\left(X - \mathbb{E}X\right)\left(Y - \mathbb{E}Y\right)\right) = \mathbb{E}XY - \mathbb{E}X\mathbb{E}Y.$$

The correlation between $X$ and $Y$ is

$$\text{corr}\left(X,Y\right) = \rho_{XY} = \frac{\sigma_{XY}}{\sigma_x\sigma_Y}.$$

The Cauchy-Schwarz Inequality implies that

$$|\rho_{XY}| \leq 1. \qquad (B.7)$$

The correlation is a measure of linear dependence, free of units of measurement.

If $X$ and $Y$ are independent, then $\sigma_{XY} = 0$ and $\rho_{XY} = 0$. The reverse, however, is not true. For example, if $\mathbb{E}X = 0$ and $\mathbb{E}X^3 = 0$, then $\mathrm{cov}(X, X^2) = 0$.

A useful fact is that

$$\mathrm{var}\,(X + Y) = \mathrm{var}(X) + \mathrm{var}(Y) + 2\,\mathrm{cov}(X, Y).$$

An implication is that if $X$ and $Y$ are independent, then

$$\mathrm{var}\,(X + Y) = \mathrm{var}(X) + \mathrm{var}(Y),$$

the variance of the sum is the sum of the variances.

A $k \times 1$ random vector $\mathbf{X} = (X_1, ..., X_k)'$ is a function from $S$ to $\mathbb{R}^k$. Let $\mathbf{x} = (x_1, ..., x_k)'$ denote a vector in $\mathbb{R}^k$. (In this Appendix, we use bold to denote vectors. Bold capitals $\mathbf{X}$ are random vectors and bold lower case $\mathbf{x}$ are nonrandom vectors. Again, this is in distinction to the notation used in the bulk of the text) The vector $\mathbf{X}$ has the distribution and density functions

$$
\begin{aligned}
F(\mathbf{x}) &= \mathbb{P}(\mathbf{X} \leq \mathbf{x}) \\
f(\mathbf{x}) &= \frac{\partial^k}{\partial x_1 \cdots \partial x_k} F(\mathbf{x}).
\end{aligned}
$$

For a measurable function $\mathbf{g} : \mathbb{R}^k \to \mathbb{R}^s$, we define the expectation

$$\mathbb{E}\mathbf{g}(\mathbf{X}) = \int_{\mathbb{R}^k} g(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

where the symbol $d\mathbf{x}$ denotes $dx_1 \cdots dx_k$. In particular, we have the $k \times 1$ multivariate mean

$$\boldsymbol{\mu} = \mathbb{E}\mathbf{X}$$

and $k \times k$ covariance matrix

$$
\begin{aligned}
\boldsymbol{\Sigma} &= \mathbb{E}\left((\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})'\right) \\
&= \mathbb{E}\mathbf{X}\mathbf{X}' - \boldsymbol{\mu}\boldsymbol{\mu}'
\end{aligned}
$$

If the elements of $\mathbf{X}$ are mutually independent, then $\boldsymbol{\Sigma}$ is a diagonal matrix and

$$\mathrm{var}\left(\sum_{i=1}^{k} \mathbf{X}_i\right) = \sum_{i=1}^{k} \mathrm{var}\,(\mathbf{X}_i)$$

## B.7 Conditional Distributions and Expectation

The **conditional density** of $Y$ given $\mathbf{X} = \mathbf{x}$ is defined as

$$f_{Y|\mathbf{X}}(y \mid \mathbf{x}) = \frac{f(\mathbf{x}, y)}{f_{\mathbf{X}}(\mathbf{x})}$$

if $f_X(x) > 0$. One way to derive this expression from the definition of conditional probability is

$$
\begin{aligned}
f_{Y|X}(y \mid x) &= \frac{\partial}{\partial y} \lim_{\varepsilon \to 0} \mathbb{P}(Y \le y \mid x \le X \le x + \varepsilon) \\
&= \frac{\partial}{\partial y} \lim_{\varepsilon \to 0} \frac{\mathbb{P}(\{Y \le y\} \cap \{x \le X \le x + \varepsilon\})}{\mathbb{P}(x \le X \le x + \varepsilon)} \\
&= \frac{\partial}{\partial y} \lim_{\varepsilon \to 0} \frac{F(x + \varepsilon, y) - F(x, y)}{F_X(x + \varepsilon) - F_X(x)} \\
&= \frac{\partial}{\partial y} \lim_{\varepsilon \to 0} \frac{\frac{\partial}{\partial x} F(x + \varepsilon, y)}{f_X(x + \varepsilon)} \\
&= \frac{\frac{\partial^2}{\partial x \partial y} F(x, y)}{f_X(x)} \\
&= \frac{f(x, y)}{f_X(x)}.
\end{aligned}
$$

The **conditional mean** or **conditional expectation** is the function

$$
m(x) = \mathbb{E}(Y \mid X = x) = \int_{-\infty}^{\infty} y f_{Y|X}(y \mid x)\, dy.
$$

The conditional mean $m(x)$ is a function, meaning that when $X$ equals $x$, then the expected value of $Y$ is $m(x)$.

Similarly, we define the conditional variance of $Y$ given $X = x$ as

$$
\begin{aligned}
\sigma^2(x) &= \operatorname{var}(Y \mid X = x) \\
&= \mathbb{E}\left((Y - m(x))^2 \mid X = x\right) \\
&= \mathbb{E}(Y^2 \mid X = x) - m(x)^2.
\end{aligned}
$$

Evaluated at $x = X$, the conditional mean $m(X)$ and conditional variance $\sigma^2(X)$ are random variables, functions of $X$. We write this as $\mathbb{E}(Y \mid X) = m(X)$ and $\operatorname{var}(Y \mid X) = \sigma^2(X)$. For example, if $\mathbb{E}(Y \mid X = x) = \alpha + \beta' x$, then $\mathbb{E}(Y \mid X) = \alpha + \beta' X$, a transformation of $X$.

The following are important facts about conditional expectations.

**Simple Law of Iterated Expectations:**

$$
\mathbb{E}(\mathbb{E}(Y \mid X)) = \mathbb{E}(Y) \tag{B.8}
$$

**Proof:**

$$
\begin{aligned}
\mathbb{E}(\mathbb{E}(Y \mid X)) &= \mathbb{E}(m(X)) \\
&= \int_{-\infty}^{\infty} m(x) f_X(x)\, dx \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{Y|X}(y \mid x) f_X(x)\, dy\, dx \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f(y, x)\, dy\, dx \\
&= \mathbb{E}(Y).
\end{aligned}
$$

**Law of Iterated Expectations:**

$$
\mathbb{E}(\mathbb{E}(Y \mid X, Z) \mid X) = \mathbb{E}(Y \mid X) \tag{B.9}
$$

210

**Conditioning Theorem.** For any function $g(\boldsymbol{x})$,

$$\mathbb{E}\left(g(\boldsymbol{X})Y \mid \boldsymbol{X}\right) = g\left(\boldsymbol{X}\right)\mathbb{E}\left(Y \mid \boldsymbol{X}\right) \tag{B.10}$$

**Proof:** Let

$$
\begin{aligned}
h(\boldsymbol{x}) &= \mathbb{E}\left(g(\boldsymbol{X})Y \mid \boldsymbol{X} = \boldsymbol{x}\right) \\
&= \int_{-\infty}^{\infty} g(\boldsymbol{x}) y f_{Y\mid X}\left(y \mid \boldsymbol{x}\right) dy \\
&= g(\boldsymbol{x}) \int_{-\infty}^{\infty} y f_{Y\mid X}\left(y \mid \boldsymbol{x}\right) dy \\
&= g(\boldsymbol{x}) m(\boldsymbol{x})
\end{aligned}
$$

where $m(\boldsymbol{x}) = \mathbb{E}\left(Y \mid \boldsymbol{X} = \boldsymbol{x}\right)$. Thus $h(\boldsymbol{X}) = g(\boldsymbol{X})m(\boldsymbol{X})$, which is the same as $\mathbb{E}\left(g(\boldsymbol{X})Y \mid \boldsymbol{X}\right) = g\left(\boldsymbol{X}\right)\mathbb{E}\left(Y \mid \boldsymbol{X}\right)$.

## B.8 Transformations

Suppose that $\boldsymbol{X} \in \mathbb{R}^k$ with continuous distribution function $F_{\boldsymbol{X}}(\boldsymbol{x})$ and density $f_{\boldsymbol{X}}(\boldsymbol{x})$. Let $\boldsymbol{Y} = \boldsymbol{g}(\boldsymbol{X})$ where $\boldsymbol{g}(\boldsymbol{x}) : \mathbb{R}^k \to \mathbb{R}^k$ is one-to-one, differentiable, and invertible. Let $\boldsymbol{h}(\boldsymbol{y})$ denote the inverse of $\boldsymbol{g}(\boldsymbol{x})$. The **Jacobian** is

$$J(\boldsymbol{y}) = \det\left(\frac{\partial}{\partial \boldsymbol{y}'}\boldsymbol{h}(\boldsymbol{y})\right).$$

Consider the univariate case $k = 1$. If $g(x)$ is an increasing function, then $g(X) \leq Y$ if and only if $X \leq h(Y)$, so the distribution function of $Y$ is

$$
\begin{aligned}
F_Y(y) &= \mathbb{P}\left(g(X) \leq y\right) \\
&= \mathbb{P}\left(X \leq h(Y)\right) \\
&= F_X\left(h(Y)\right).
\end{aligned}
$$

Taking the derivative, the density of $Y$ is

$$f_Y(y) = \frac{d}{dy}F_Y(y) = f_X\left(h(Y)\right)\frac{d}{dy}h(y).$$

If $g(x)$ is a decreasing function, then $g(X) \leq Y$ if and only if $X \geq h(Y)$, so

$$
\begin{aligned}
F_Y(y) &= \mathbb{P}\left(g(X) \leq y\right) \\
&= 1 - \mathbb{P}\left(X \geq h(Y)\right) \\
&= 1 - F_X\left(h(Y)\right)
\end{aligned}
$$

and the density of $Y$ is

$$f_Y(y) = -f_X\left(h(Y)\right)\frac{d}{dy}h(y).$$

We can write these two cases jointly as

$$f_Y(y) = f_X\left(h(Y)\right)\left|J(y)\right|. \tag{B.11}$$

This is known as the **change-of-variables** formula. This same formula (B.11) holds for $k > 1$, but its justification requires deeper results from analysis.

As one example, take the case $X \sim U[0, 1]$ and $Y = -\log(X)$. Here, $g(x) = -\log(x)$ and $h(y) = \exp(-y)$ so the Jacobian is $J(y) = -\exp(y)$. As the range of $X$ is $[0, 1]$, that for $Y$ is $[0, \infty)$. Since $f_X(x) = 1$ for $0 \leq x \leq 1$ (B.11) shows that

$$f_Y(y) = \exp(-y), \qquad 0 \leq y \leq \infty,$$

an exponential density.

## B.9   Normal and Related Distributions

The **standard normal** density is

$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right), \qquad -\infty < x < \infty.$$

It is conventional to write $X \sim N(0,1)$, and to denote the standard normal density function by $\phi(x)$ and its distribution function by $\Phi(x)$. The latter has no closed-form solution. The normal density has all moments finite. Since it is symmetric about zero all odd moments are zero. By iterated integration by parts, we can also show that $\mathbb{E}X^2 = 1$ and $\mathbb{E}X^4 = 3$. In fact, for any positive integer $m$, $\mathbb{E}X^{2m} = (2m-1)!! = (2m-1) \cdot (2m-3) \cdots 1$. Thus $\mathbb{E}X^4 = 3$, $\mathbb{E}X^6 = 15$, $\mathbb{E}X^8 = 105$, and $\mathbb{E}X^{10} = 945$.

If $Z$ is standard normal and $X = \mu + \sigma Z$, then using the change-of-variables formula, $X$ has density

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \qquad -\infty < x < \infty.$$

which is the **univariate normal density**. The mean and variance of the distribution are $\mu$ and $\sigma^2$, and it is conventional to write $X \sim N(\mu, \sigma^2)$.

For $x \in \mathbb{R}^k$, the **multivariate normal density** is

$$f(x) = \frac{1}{(2\pi)^{k/2} \det(\Sigma)^{1/2}} \exp\left(-\frac{(x-\mu)' \Sigma^{-1} (x-\mu)}{2}\right), \qquad x \in \mathbb{R}^k.$$

The mean and covariance matrix of the distribution are $\mu$ and $\Sigma$, and it is conventional to write $X \sim N(\mu, \Sigma)$.

The MGF and CF of the multivariate normal are $\exp\left(\lambda'\mu + \lambda'\Sigma\lambda/2\right)$ and $\exp\left(i\lambda'\mu - \lambda'\Sigma\lambda/2\right)$, respectively.

If $X \in \mathbb{R}^k$ is multivariate normal and the elements of $X$ are mutually uncorrelated, then $\Sigma = \text{diag}\{\sigma_j^2\}$ is a diagonal matrix. In this case the density function can be written as

$$\begin{aligned}
f(x) &= \frac{1}{(2\pi)^{k/2} \sigma_1 \cdots \sigma_k} \exp\left(-\left(\frac{(x_1-\mu_1)^2/\sigma_1^2 + \cdots + (x_k-\mu_k)^2/\sigma_k^2}{2}\right)\right) \\
&= \prod_{j=1}^{k} \frac{1}{(2\pi)^{1/2} \sigma_j} \exp\left(-\frac{(x_j-\mu_j)^2}{2\sigma_j^2}\right)
\end{aligned}$$

which is the product of marginal univariate normal densities. This shows that if $X$ is multivariate normal with uncorrelated elements, then they are mutually independent.

**Theorem B.9.1** *If $X \sim N(\mu, \Sigma)$ and $Y = a + BX$ with $B$ an invertible matrix, then $Y \sim N(a + B\mu, B\Sigma B')$.*

**Theorem B.9.2** *Let $X \sim N(0, I_r)$. Then $Q = X'X$ is distributed chi-square with $r$ degrees of freedom, written $\chi_r^2$.*

**Theorem B.9.3** *If $Z \sim N(0, A)$ with $A > 0$, $q \times q$, then $Z'A^{-1}Z \sim \chi_q^2$.*

**Theorem B.9.4** *Let $Z \sim N(0,1)$ and $Q \sim \chi_r^2$ be independent. Then $T_r = Z/\sqrt{Q/r}$ is distributed as student's t with $r$ degrees of freedom.*

**Proof of Theorem B.9.1.** By the change-of-variables formula, the density of $Y = a + BX$ is

$$f(y) = \frac{1}{(2\pi)^{k/2} \det (\Sigma_Y)^{1/2}} \exp \left( -\frac{(y - \mu_Y)' \Sigma_Y^{-1} (y - \mu_Y)}{2} \right), \qquad y \in \mathbb{R}^k.$$

where $\mu_Y = a + B\mu$ and $\Sigma_Y = B\Sigma B'$, where we used the fact that $\det (B\Sigma B')^{1/2} = \det (\Sigma)^{1/2} \det (B)$. ∎

**Proof of Theorem B.9.2.** First, suppose a random variable $Q$ is distributed chi-square with $r$ degrees of freedom. It has the MGF

$$\mathbb{E} \exp (tQ) = \int_0^\infty \frac{1}{\Gamma \left(\frac{r}{2}\right) 2^{r/2}} x^{r/2-1} \exp (tx) \exp (-x/2) \, dy = (1 - 2t)^{-r/2}$$

where the second equality uses the fact that $\int_0^\infty y^{a-1} \exp (-by) \, dy = b^{-a} \Gamma(a)$, which can be found by applying change-of-variables to the gamma function. Our goal is to calculate the MGF of $Q = X'X$ and show that it equals $(1 - 2t)^{-r/2}$, which will establish that $Q \sim \chi_r^2$.

Note that we can write $Q = X'X = \sum_{j=1}^r Z_j^2$ where the $Z_j$ are independent $N(0,1)$. The distribution of each of the $Z_j^2$ is

$$\begin{aligned} \mathbb{P} \left( Z_j^2 \le y \right) &= 2\mathbb{P} \left( 0 \le Z_j \le \sqrt{y} \right) \\ &= 2 \int_0^{\sqrt{y}} \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{x^2}{2} \right) dx \\ &= \int_0^y \frac{1}{\Gamma \left(\frac{1}{2}\right) 2^{1/2}} s^{-1/2} \exp \left( -\frac{s}{2} \right) ds \end{aligned}$$

using the change–of–variables $s = x^2$ and the fact $\Gamma \left(\frac{1}{2}\right) = \sqrt{\pi}$. Thus the density of $Z_j^2$ is

$$f_1(x) = \frac{1}{\Gamma \left(\frac{1}{2}\right) 2^{1/2}} x^{-1/2} \exp \left( -\frac{x}{2} \right)$$

which is the $\chi_1^2$ and by our above calculation has the MGF of $\mathbb{E} \exp \left( tZ_j^2 \right) = (1 - 2t)^{-1/2}$.

Since the $Z_j^2$ are mutually independent, (B.6) implies that the MGF of $Q = \sum_{j=1}^r Z_j^2$ is $\left[ (1 - 2t)^{-1/2} \right]^r = (1 - 2t)^{-r/2}$, which is the MGF of the $\chi_r^2$ density as desired. ∎

**Proof of Theorem B.9.3.** The fact that $A > 0$ means that we can write $A = CC'$ where $C$ is non-singular. Then $A^{-1} = C^{-1'}C^{-1}$ and

$$C^{-1}Z \sim N \left( 0, C^{-1}AC^{-1'} \right) = N \left( 0, C^{-1}CC'C^{-1'} \right) = N \left( 0, I_q \right).$$

Thus

$$Z'A^{-1}Z = Z'C^{-1'}C^{-1}Z = \left( C^{-1}Z \right)' \left( C^{-1}Z \right) \sim \chi_q^2.$$

∎

**Proof of Theorem B.9.4.** Using the simple law of iterated expectations, $T_r$ has distribution

function

$$
\begin{aligned}
F\left(x\right) &= \mathbb{P}\left(\frac{Z}{\sqrt{Q/r}} \le x\right) \\
&= \mathbb{E}\left\{Z \le x\sqrt{\frac{Q}{r}}\right\} \\
&= \mathbb{E}\left[\mathbb{P}\left(Z \le x\sqrt{\frac{Q}{r}} \mid Q\right)\right] \\
&= \mathbb{E}\Phi\left(x\sqrt{\frac{Q}{r}}\right)
\end{aligned}
$$

Thus its density is

$$
\begin{aligned}
f\left(x\right) &= \mathbb{E}\frac{d}{dx}\Phi\left(x\sqrt{\frac{Q}{r}}\right) \\
&= \mathbb{E}\left(\phi\left(x\sqrt{\frac{Q}{r}}\right)\sqrt{\frac{Q}{r}}\right) \\
&= \int_0^\infty \left(\frac{1}{\sqrt{2\pi}}\exp\left(-\frac{qx^2}{2r}\right)\right)\sqrt{\frac{q}{r}}\left(\frac{1}{\Gamma\left(\frac{r}{2}\right)2^{r/2}}q^{r/2-1}\exp\left(-q/2\right)\right)dq \\
&= \frac{\Gamma\left(\frac{r+1}{2}\right)}{\sqrt{r\pi}\Gamma\left(\frac{r}{2}\right)}\left(1+\frac{x^2}{r}\right)^{-\left(\frac{r+1}{2}\right)}
\end{aligned}
$$

which is that of the student t with $r$ degrees of freedom. ∎

214