

Chapter 1

Probability, Random Variables and Expectations

Note: The primary reference for these notes is Mittelhammer (1999). Other treatments of probability theory include Gallant (1997), Casella and Berger (2001) and Grimmett and Stirzaker (2001).

This chapter provides an overview of probability theory as it applied to both discrete and continuous random variables. The material covered in this chapter serves as a foundation of the econometric sequence and is useful throughout financial economics. The chapter begins with a discussion of the axiomatic foundations of probability theory and then proceeds to describe properties of univariate random variables. Attention then turns to multivariate random variables and important difference from univariate random variables. Finally, the chapter discusses the expectations operator and moments.

1.1 Axiomatic Probability

Probability theory is derived from a small set of axioms – a minimal set of essential assumptions. A deep understanding of axiomatic probability theory is *not* essential to financial econometrics or to the use of probability and statistics in general, although understanding these core concepts does provide additional insight.

The first concept in probability theory is the sample space, which is an abstract concept containing primitive probability events.

Definition 1.1 (Sample Space). The sample space is a set, Ω , that contains all possible outcomes.

Example 1.1. Suppose interest is on a standard 6-sided die. The sample space is 1-dot, 2-dots, ..., 6-dots.

Example 1.2. Suppose interest is in a standard 52-card deck. The sample space is then $A_{\clubsuit}, 2_{\clubsuit}, 3_{\clubsuit}, \dots, J_{\clubsuit}, Q_{\clubsuit}, K_{\clubsuit}, A_{\diamond}, \dots, K_{\diamond}, A_{\heartsuit}, \dots, K_{\heartsuit}, A_{\spadesuit}, \dots, K_{\spadesuit}$.

Example 1.3. Suppose interest is in the logarithmic stock return, defined as $r_t = \ln P_t - \ln P_{t-1}$, then the sample space is \mathbb{R} , the real line.

The next item of interest is an event.

Definition 1.2 (Event). An event, ω , is a subset of the sample space Ω .

An event may be any subsets of the sample space Ω (including the entire sample space), and the set of all events is known as the event space.

Definition 1.3 (Event Space). The set of all events in the sample space Ω is called the event space, and is denoted \mathcal{F} .

Event spaces are a somewhat more difficult concept. For finite event spaces, the event space is usually the power set of the outcomes – that is, the set of all possible unique sets that can be constructed from the elements. When variables can take infinitely many outcomes, then a more nuanced definition is needed, although the main idea is to define the event space to be all non-empty intervals (so that each interval has infinitely many points in it).

Example 1.4. Suppose interest lies in the outcome of a coin flip. Then the sample space is $\{H, T\}$ and the event space is $\{\emptyset, \{H\}, \{T\}, \{H, T\}\}$ where \emptyset is the empty set.

The first two axioms of probability are simple: all probabilities must be non-negative and the total probability of all events is one.

Axiom 1.1. For any event $\omega \in \mathcal{F}$,

$$\Pr(\omega) \geq 0. \quad (1.1)$$

Axiom 1.2. The probability of all events in the sample space Ω is unity, i.e.

$$\Pr(\Omega) = 1. \quad (1.2)$$

The second axiom is a normalization that states that the probability of the entire sample space is 1 and ensures that the sample space must contain all events that may occur. $\Pr(\cdot)$ is a set-valued function – that is, $\Pr(\omega)$ returns the probability, a number between 0 and 1, of observing an event ω .

Before proceeding, it is useful to refresh four concepts from set theory.

Definition 1.4 (Set Union). Let A and B be two sets, then the union is defined

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

A union of two sets contains all elements that are in either set.

Definition 1.5 (Set Intersection). Let A and B be two sets, then the intersection is defined

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

The intersection contains only the elements that are in both sets.

Definition 1.6 (Set Complement). Let A be a set, then the complement set, denoted

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

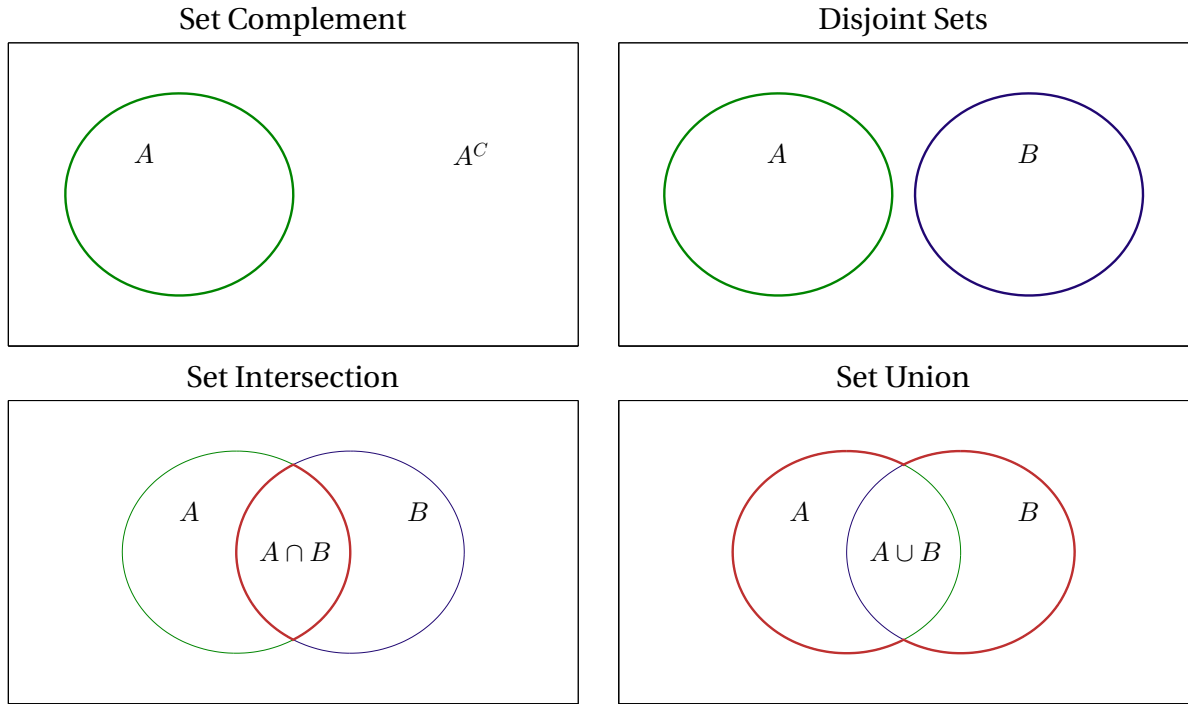


Figure 1.1: The four set definitions presented in \mathbb{R}^2 . The upper left panel shows a set and its complement. The upper right shows two disjoint sets. The lower left shows the intersection of two sets (darkened region) and the lower right shows the union of two sets (darkened region). In all diagrams, the outer box represents the entire space.

The complement of a set contains all elements which are not contained in the set.

Definition 1.7 (Disjoint Sets). Let A and B be sets, then A and B are disjoint if and only if $A \cap B = \emptyset$.

Figure 1.1 provides a graphical representation of the four set operations in a 2-dimensional space.

The third and final axiom states that probability is additive when sets are disjoint.

Axiom 1.3. Let $\{A_i\}$, $i = 1, 2, \dots$ be a finite or countably infinite set of disjoint events.¹ Then

$$\Pr \left(\bigcup_{i=1}^{\infty} A_i \right) = \sum_{i=1}^{\infty} \Pr(A_i). \quad (1.3)$$

Assembling a sample space, event space and a probability measure into a set produces what is known as a probability space. Throughout the course, and in virtually all statistics, a complete

1

Definition 1.8. A S set is countably infinite if there exists a bijective (one-to-one) function from the elements of S to the natural numbers $\mathbb{N} = \{1, 2, \dots\}$. Common sets that are countable infinite include the integers (\mathbb{Z}) and the rational numbers (\mathbb{Q}).

probability space is assumed (typically without explicitly stating this assumption).²

Definition 1.9 (Probability Space). A probability space is denoted using the tuple $(\Omega, \mathcal{F}, \Pr)$ where Ω is the sample space, \mathcal{F} is the event space and \Pr is the probability set function which has domain $\omega \in \mathcal{F}$.

The three axioms of modern probability are very powerful, and a large number of theorems can be proven using only these axioms. A few simple examples are provided, and selected proofs appear in the Appendix.

Theorem 1.1. *Let A be an event in the sample space Ω , and let A^c be the complement of A so that $\Omega = A \cup A^c$. Then $\Pr(A) = 1 - \Pr(A^c)$.*

Since A and A^c are disjoint, and by definition A^c is everything not in A , then the probability of the two must be unity.

Theorem 1.2. *Let A and B be events in the sample space Ω . Then $\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$.*

This theorem shows that for any two sets, the probability of the union of the two sets is equal to the probability of the two sets minus the probability of the intersection of the sets.

1.1.1 Conditional Probability

Conditional probability extends the basic concepts of probability to the case where interest lies in the probability of one event conditional on the occurrence of another event.

Definition 1.10 (Conditional Probability). Let A and B be two events in the sample space Ω . If $\Pr(B) \neq 0$, then the conditional probability of the event A , given event B , is given by

$$\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)}. \quad (1.4)$$

The definition of conditional probability is intuitive. The probability of observing an event in set A , given an event in the set B has occurred, is the probability of observing an event in the intersection of the two sets normalized by the probability of observing an event in set B .

Example 1.5. In the example of rolling a die, suppose $A = \{1, 3, 5\}$ is the event that the outcome is odd and $B = \{1, 2, 3\}$ is the event that the outcome of the roll is less than 4. Then the conditional probability of A given B is

$$\frac{\Pr(\{1, 3\})}{\Pr(\{1, 2, 3\})} = \frac{\frac{2}{6}}{\frac{3}{6}} = \frac{2}{3}$$

since the intersection of A and B is $\{1, 3\}$.

The axioms can be restated in terms of conditional probability, where the sample space consists of the events in the set B .

²A complete probability space is complete if and only if $B \in \mathcal{F}$ where $\Pr(B) = 0$ and $A \subset B$, then $A \in \mathcal{F}$. This condition ensures that probability can be assigned to any event.

1.1.2 Independence

Independence of two measurable sets means that any information about an event occurring in one set has no information about whether an event occurs in another set.

Definition 1.11. Let A and B be two events in the sample space Ω . Then A and B are independent if and only if

$$\Pr(A \cap B) = \Pr(A) \Pr(B) \quad (1.5)$$

, $A \perp\!\!\!\perp B$ is commonly used to indicate that A and B are independent.

One immediate implication of the definition of independence is that when A and B are independent, then the conditional probability of one given the other is the same as the unconditional probability of the random variable – i.e. $\Pr(A|B) = \Pr(A)$.

1.1.3 Bayes Rule

Bayes rule is frequently encountered in both statistics (known as Bayesian statistics) and in financial models where agents learn about their environment. Bayes rule follows as a corollary to a theorem that states that the total probability of a set A is equal to the conditional probability of A given a set of disjoint sets B which span the sample space.

Theorem 1.3. Let $B_i, i = 1, 2, \dots$ be a finite or countably infinite partition of the sample space Ω so that $B_j \cap B_k = \emptyset$ for $j \neq k$ and $\bigcup_{i=1}^{\infty} B_i = \Omega$. Let $\Pr(B_i) > 0$ for all i , then for any set A ,

$$\Pr(A) = \sum_{i=1}^{\infty} \Pr(A|B_i) \Pr(B_i). \quad (1.6)$$

Bayes rule restates the previous theorem so that the probability of observing an event in B_j given an event in A is observed can be related to the conditional probability of A given B_j .

Corollary 1.1 (Bayes Rule). Let $B_i, i = 1, 2, \dots$ be a finite or countably infinite partition of the sample space Ω so that $B_j \cap B_k = \emptyset$ for $j \neq k$ and $\bigcup_{i=1}^{\infty} B_i = \Omega$. Let $\Pr(B_i) > 0$ for all i , then for any set A where $\Pr(A) > 0$,

$$\begin{aligned} \Pr(B_j|A) &= \frac{\Pr(A|B_j) \Pr(B_j)}{\sum_{i=1}^{\infty} \Pr(A|B_i) \Pr(B_i)} \\ &= \frac{\Pr(A|B_j) \Pr(B_j)}{\Pr(A)} \end{aligned}$$

An immediate consequence of the definition of conditional probability is the

$$\Pr(A \cap B) = \Pr(A|B) \Pr(B),$$

which is referred to as the multiplication rule. Also notice that the order of the two sets is arbitrary, so that the rule can be equivalently stated as $\Pr(A \cap B) = \Pr(B|A) \Pr(A)$. Combining these two (as long as $\Pr(A) > 0$),

$$\begin{aligned} \Pr(A|B) \Pr(B) &= \Pr(B|A) \Pr(A) \\ \Rightarrow \Pr(B|A) &= \frac{\Pr(A|B) \Pr(B)}{\Pr(A)}. \end{aligned} \quad (1.7)$$

Example 1.6. Suppose a family has 2 children and one is a boy, and that the probability of having a child of either sex is equal and independent across children. What is the probability that they have 2 boys?

Before learning that one child is a boy, there are 4 equally probable possibilities: $\{B, B\}$, $\{B, G\}$, $\{G, B\}$ and $\{G, G\}$. Using Bayes rule,

$$\begin{aligned} \Pr(\{B, B\} | B \geq 1) &= \frac{\Pr(B \geq 1 | \{B, B\}) \times \Pr(\{B, B\})}{\sum_{S \in \{\{B, B\}, \{B, G\}, \{G, B\}, \{G, G\}\}} \Pr(B \geq 1 | S) \Pr(S)} \\ &= \frac{1 \times \frac{1}{4}}{1 \times \frac{1}{4} + 1 \times \frac{1}{4} + 1 \times \frac{1}{4} + 0 \times \frac{1}{4}} \\ &= \frac{1}{3} \end{aligned}$$

so that knowing one child is a boy increases the probability of 2 boys from $\frac{1}{4}$ to $\frac{1}{3}$. Note that

$$\sum_{S \in \{\{B, B\}, \{B, G\}, \{G, B\}, \{G, G\}\}} \Pr(B \geq 1 | S) \Pr(S) = \Pr(B \geq 1).$$

Example 1.7. The famous Monte Hall *Let's Make a Deal* television program is an example of Bayes rule. Contestants competed for one of three prizes, a large one (e.g. a car) and two uninteresting ones (duds). The prizes were hidden behind doors numbered 1, 2 and 3. Before the contest starts, the contestant has no information about the which door has the large prize, and to the initial probabilities are all $\frac{1}{3}$. During the negotiations with the host, it is revealed that one of the non-selected doors does *not* contain the large prize. The host then gives the contestant the chance to switch from the door initially chosen to the one remaining door. For example, suppose the contestant choose door 1 initially, and that the host revealed that the large prize is not behind door 3. The contestant then has the chance to choose door 2 or to stay with door 1. In this example, B is the event where the contestant chooses the door which hides the large prize, and A is the event that the large prize is not behind door 2.

Initially there are three equally likely outcomes (from the contestant's point of view), where D indicates dud, L indicates the large prize, and the order corresponds to the door number.

$$\{D, D, L\}, \{D, L, D\}, \{L, D, D\}$$

The contestant has a $\frac{1}{3}$ chance of having the large prize behind door 1. The host will never remove the large prize, and so applying Bayes rule we have

$$\begin{aligned} \Pr(L = 2 | H = 3, S = 1) &= \frac{\Pr(H = 3 | S = 1, L = 2) \times \Pr(L = 2 | S = 1)}{\sum_{i=1}^3 \Pr(H = 3 | S = 1, L = i) \times \Pr(L = i | S = 1)} \\ &= \frac{1 \times \frac{1}{3}}{\frac{1}{2} \times \frac{1}{3} + 1 \times \frac{1}{3} + 0 \times \frac{1}{3}} \\ &= \frac{\frac{1}{3}}{\frac{1}{2}} \\ &= \frac{2}{3}. \end{aligned}$$

where H is the door the host reveals, S is initial door selected, and L is the door containing the large prize. This shows that the probability the large prize is behind door 2, given that the player initially selected door 1 and the host revealed door 3 can be computed using Bayes rule.

$\Pr(H = 3|S = 1, L = 2)$ is the probability that the host shows door 3 given the contestant selected door 1 and the large prize is behind door 2, which always happens since the host will never reveal the large prize. $P(L = 2|S = 1)$ is the probability that the large is in door 2 given the contestant selected door 1, which is $\frac{1}{3}$. $\Pr(H = 3|S = 1, L = 1)$ is the probability that the host reveals door 3 given that door 1 was selected and contained the large prize, which is $\frac{1}{2}$, and $P(H = 3|S = 1, L = 3)$ is the probability that the host reveals door 3 given door 3 contains the prize, which never happens.

Bayes rule shows that it is always optimal to switch doors. This is a counter-intuitive result and occurs since the host's action reveals information about the location of the large prize. Essentially, the two doors not selected by the host have combined probability $\frac{2}{3}$ of containing the large prize before the doors are opened – opening the third assigns its probability to the door not opened.

1.2 Univariate Random Variables

Studying the behavior of random variables, and more importantly functions of random variables (i.e. statistics) is essential for both the theory and practice of financial econometrics. This section covers univariate random variables and multivariate random variables are discussed later.

The previous discussion of probability is set based and so includes objects which cannot be described as random variables, which are a limited (but highly useful) sub-class of all objects that can be described using probability theory. The primary characteristic of a random variable is that it takes values on the real line.

Definition 1.12 (Random Variable). Let $(\Omega, \mathcal{F}, \Pr)$ be a probability space. If $X : \Omega \rightarrow \mathbb{R}$ is a real-valued function have as its domain elements of Ω , then X is called a random variable.

A random variable is essentially a function which takes $\omega \in \Omega$ as an input and returns a value $x \in \mathbb{R}$, where \mathbb{R} is the symbol for the real line. Random variables come in one of three forms: discrete, continuous and mixed. Random variables which mix discrete and continuous distributions are generally less important in financial economics and so here the focus is on discrete and continuous random variables.

Definition 1.13 (Discrete Random Variable). A random variable is called discrete if its range consists of a countable (possibly infinite) number of elements.

While discrete random variables are less useful than continuous random variables, they are still commonly encountered.

Example 1.8. A random variable which takes on values in $\{0, 1\}$ is known as a Bernoulli random variable, and is the simplest non-degenerate random variable (see Section 1.2.3.1).³ Bernoulli random variables are often used to model “success” or “failure”, where success is loosely defined – a large negative return, the existence of a bull market or a corporate default.

³A degenerate random variable always takes the same value, and so is not meaningfully random.

The distinguishing characteristic of a discrete random variable is not that it takes only finitely many values, but that the values it takes are distinct in the sense that it is possible to fit small intervals around each point without the overlap.

Example 1.9. Poisson random variables take values in $\{0, 1, 2, 3, \dots\}$ (an infinite range), and are commonly used to model hazard rates (i.e. the number of occurrences of an event in an interval). They are especially useful in modeling trading activity (see Section 1.2.3.2).

1.2.1 Mass, Density, and Distribution Functions

Discrete random variables are characterized by a probability mass function (pmf) which gives the probability of observing a particular value of the random variable.

Definition 1.14 (Probability Mass Function). The probability mass function, f , for a discrete random variable X is defined as $f(x) = \Pr(x)$ for all $x \in R(X)$, and $f(x) = 0$ for all $x \notin R(X)$ where $R(X)$ is the range of X (i.e. the values for which X is defined).

Example 1.10. The probability mass function of a Bernoulli random variable takes the form

$$f(x; p) = p^x (1 - p)^{1-x}$$

where $p \in [0, 1]$ is the probability of success.

Figure 1.2 contains a few examples of Bernoulli pmfs using data from the FTSE 100 and S&P 500 over the period 1984–2012. Both weekly returns, using Friday to Friday prices and monthly returns, using end-of-month prices, were constructed. Log returns were used ($r_t = \ln(P_t/P_{t-1})$) in both examples. Two of the pmfs defined success as the return being positive. The other two define the probability of success as a return larger than -1% (weekly) or larger than -4% (monthly). These show that the probability of a positive return is much larger for monthly horizons than for weekly.

Example 1.11. The probability mass function of a Poisson random variable is

$$f(x; \lambda) = \frac{\lambda^x}{x!} \exp(-\lambda)$$

where $\lambda \in [0, \infty)$ determines the intensity of arrival (the average value of the random variable).

The pmf of the Poisson distribution can be evaluated for every value of $x \geq 0$, which is the support of a Poisson random variable. Figure 1.4 shows empirical distribution tabulated using a histogram for the time elapsed where .1% of the daily volume traded in the S&P 500 tracking ETF SPY on May 31, 2012. This data series is a good candidate for modeling using a Poisson distribution.

Continuous random variables, on the other hand, take a continuum of values – technically an uncountable infinity of values.

Definition 1.15 (Continuous Random Variable). A random variable is called continuous if its range is uncountably infinite and there exists a non-negative-valued function $f(x)$ defined on all $x \in (-\infty, \infty)$ such that for any event $B \subset R(X)$, $\Pr(B) = \int_{x \in B} f(x) dx$ and $f(x) = 0$ for all $x \notin R(X)$ where $R(X)$ is the range of X (i.e. the values for which X is defined).

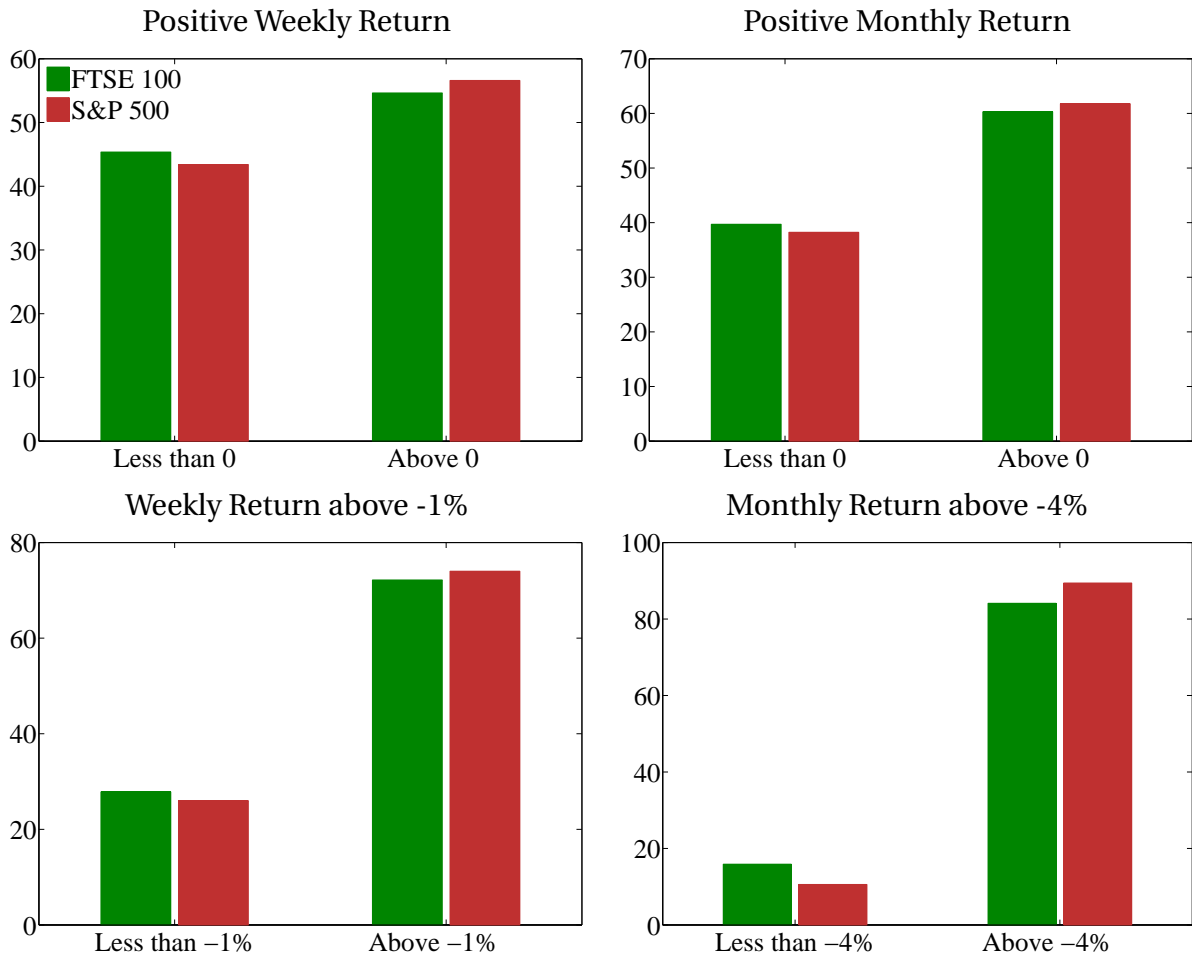


Figure 1.2: These four charts show examples of Bernoulli random variables using returns on the FTSE 100 and S&P 500. In the top two, a success was defined as a positive return. In the bottom two, a success was a return above -1% (weekly) or -4% (monthly).

The pmf of a discrete random variable is replaced with the probability density function (pdf) for continuous random variables. This change in naming reflects that the probability of a single point of a continuous random variable is 0, although the probability of observing a value inside an arbitrarily small interval in $R(X)$ is not.

Definition 1.16 (Probability Density Function). For a continuous random variable, the function f is called the probability density function (pdf).

Before providing some examples of pdfs, it is useful to characterize the properties that any pdf should have.

Definition 1.17 (Continuous Density Function Characterization). A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is a member of the class of continuous density functions if and only if $f(x) \geq 0$ for all $x \in (-\infty, \infty)$ and $\int_{-\infty}^{\infty} f(x) dx = 1$.

There are two essential properties. First, that the function is non-negative, which follows from the axiomatic definition of probability, and second, that the function integrates to 1, so that

the total probability across $R(X)$ is 1. This may seem like a limitation, but it is only a normalization since any non-negative integrable function can always be normalized to that it integrates to 1.

Example 1.12. A simple continuous random variable can be defined on $[0, 1]$ using the probability density function

$$f(x) = 12 \left(x - \frac{1}{2} \right)^2$$

and figure 1.3 contains a plot of the pdf.

This simple pdf has peaks near 0 and 1 and a trough at $1/2$. More realistic pdfs allow for values in $(-\infty, \infty)$, such as in the density of a normal random variable.

Example 1.13. The pdf of a normal random variable with parameters μ and σ^2 is given by

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(x - \mu)^2}{2\sigma^2} \right). \quad (1.8)$$

$N(\mu, \sigma^2)$ is used as a shorthand notation for a random variable with this pdf. When $\mu = 0$ and $\sigma^2 = 1$, the distribution is known as a standard normal. Figure 1.3 contains a plot of the standard normal pdf along with two other parameterizations.

For large values of x (in the absolute sense), the pdf of a standard normal takes very small values, and peaks at $x = 0$ with a value of 0.3989. The shape of the normal distribution is that of a bell (and is occasionally referred to a bell curve).

A closely related function to the pdf is the cumulative distribution function, which returns the total probability of observing a value of the random variable *less* than its input.

Definition 1.18 (Cumulative Distribution Function). The cumulative distribution function (cdf) for a random variable X is defined as $F(c) = \Pr(x \leq c)$ for all $c \in (-\infty, \infty)$.

Cumulative distribution function is used for both discrete and continuous random variables.

Definition 1.19 (Discrete cdf). When X is a discrete random variable, the cdf is

$$F(x) = \sum_{s \leq x} f(s) \quad (1.9)$$

for $x \in (-\infty, \infty)$.

Example 1.14. The cdf of a Bernoulli is

$$F(x; p) = \begin{cases} 0 & \text{if } x < 0 \\ p & \text{if } 0 \leq x < 1 \\ 1 & \text{if } x \geq 1 \end{cases}.$$

The Bernoulli cdf is simple since it only takes 3 values. The cdf of a Poisson random variable is relatively simple since it is defined as sum the probability mass function for all values less than or equal to the function's argument.

Example 1.15. The cdf of a $\text{Poisson}(\lambda)$ random variable is given by

$$F(x; \lambda) = \exp(-\lambda) \sum_{i=0}^{\lfloor x \rfloor} \frac{\lambda^i}{i!}, \quad x \geq 0.$$

where $\lfloor \cdot \rfloor$ returns the largest integer smaller than the input (the floor operator).

Continuous cdfs operate much like discrete cdfs, only the summation is replaced by an integral since there are a continuum of values possible for X .

Definition 1.20 (Continuous cdf). When X is a continuous random variable, the cdf is

$$F(x) = \int_{-\infty}^x f(s) \, ds \quad (1.10)$$

for $x \in (-\infty, \infty)$.

The integral computes the total area under the pdf starting from $-\infty$ up to x .

Example 1.16. The cdf of the random variable with pdf given by $12(x - 1/2)^2$ is

$$F(x) = 4x^3 - 6x^2 + 3x.$$

and figure 1.3 contains a plot of this cdf.

This cdf is the integral of the pdf, and checking shows that $F(0) = 0$, $F(1/2) = 1/2$ (since it is symmetric around $1/2$) and $F(1) = 1$, which must be 1 since the random variable is only defined on $[0, 1]$.^h

Example 1.17. The cdf of a normally distributed random variable with parameters μ and σ^2 is given by

$$F(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^x \exp\left(-\frac{(s - \mu)^2}{2\sigma^2}\right) \, ds. \quad (1.11)$$

Figure 1.3 contains a plot of the standard normal cdf along with two other parameterizations.

In the case of a standard normal random variable, the cdf is not available in closed form, and so when computed using a computer (i.e. in Excel or MATLAB), fast, accurate numeric approximations based on polynomial expansions are used (Abramowitz and Stegun, 1964).

The cdf can be similarly derived from the pdf as long as the cdf is continuously differentiable. At points where the cdf is not continuously differentiable, the pdf is defined to take the value 0.⁴

Theorem 1.4 (Relationship between cdf and pdf). *Let $f(x)$ and $F(x)$ represent the pdf and cdf of a continuous random variable X , respectively. The density function for X can be defined as $f(x) = \frac{\partial F(x)}{\partial x}$ whenever $f(x)$ is continuous and $f(x) = 0$ elsewhere.*

⁴Formally a pdf does not have to exist for a random variable, although a cdf always does. In practice, this is a technical point and distributions which have this property are rarely encountered in financial economics.

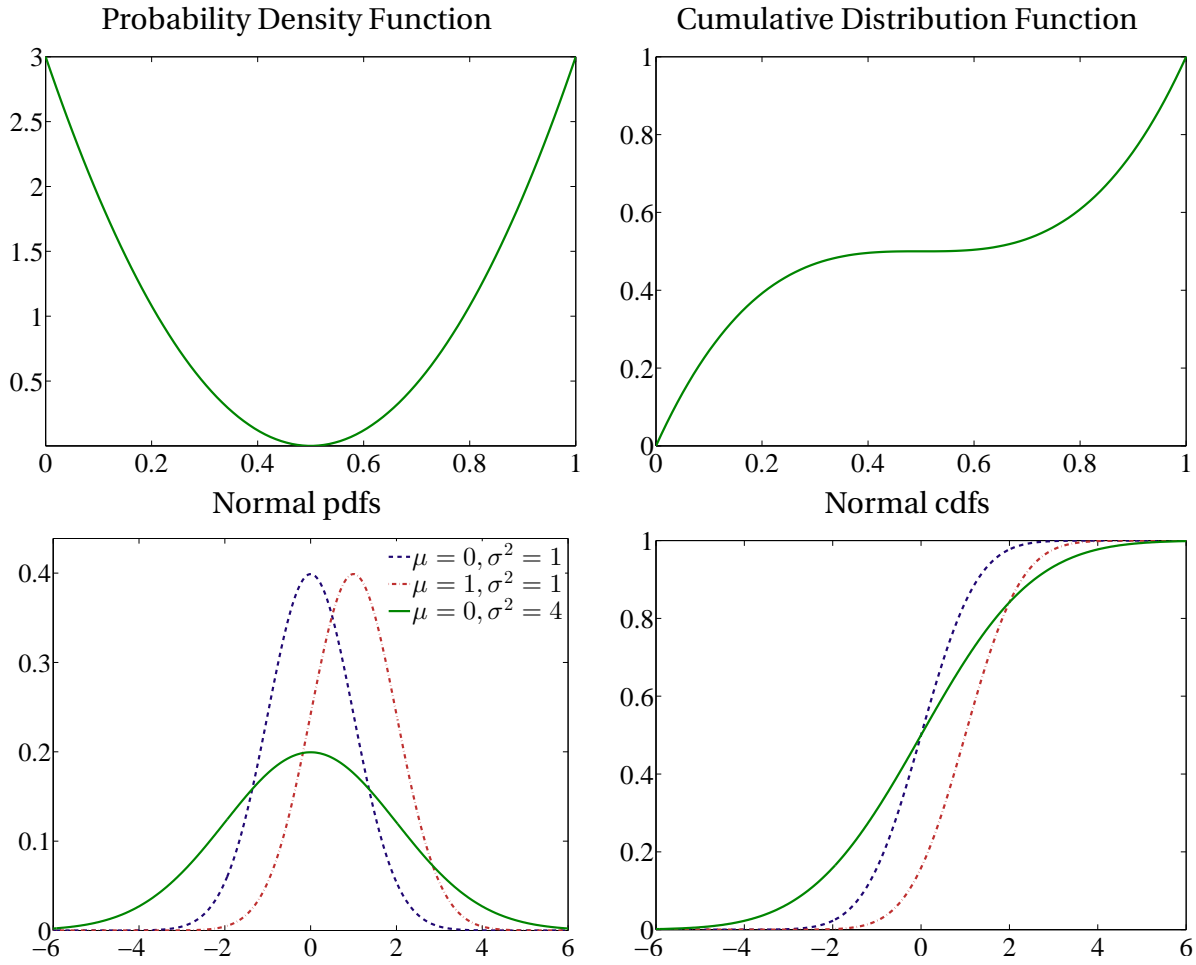


Figure 1.3: The top panels show the pdf for the density $f(x) = 12\left(x - \frac{1}{2}\right)^2$ and its associated cdf. The bottom left panel shows the probability density function for normal distributions with alternative values for μ and σ^2 . The bottom right panel shows the cdf for the same parameterizations.

Example 1.18. Taking the derivative of the cdf in the running example,

$$\begin{aligned}
 \frac{\partial F(x)}{\partial x} &= 12x^2 - 12x + 3 \\
 &= 12\left(x^2 - x + \frac{1}{4}\right) \\
 &= 12\left(x - \frac{1}{2}\right)^2.
 \end{aligned}$$

1.2.2 Quantile Functions

The quantile function is closely related to the cdf – and in many important cases, the quantile function is the inverse (function) of the cdf. Before defining quantile functions, it is necessary to

define a quantile.

Definition 1.21 (Quantile). Any number q satisfying $\Pr(x \leq q) = \alpha$ and $\Pr(x \geq q) = 1 - \alpha$ is known as the α -quantile of X and is denoted q_α .

A quantile is just the point on the cdf where the total probability that a random variable is smaller is α and the probability that the random variable takes a larger value is $1 - \alpha$. The definition of a quantile does not necessarily require uniqueness and non-unique quantiles are encountered when pdfs have regions of 0 probability (or equivalently cdfs are discontinuous). Quantiles are unique for random variables which have continuously differentiable cdfs. One common modification of the quantile definition is to select the *smallest* number which satisfies the two conditions to impose uniqueness of the quantile.

The function which returns the quantile is known as the quantile function.

Definition 1.22 (Quantile Function). Let X be a continuous random variable with cdf $F(x)$. The quantile function for X is defined as $G(\alpha) = q$ where $\Pr(x \leq q) = \alpha$ and $\Pr(x > q) = 1 - \alpha$. When $F(x)$ is one-to-one (and hence X is strictly continuous) then $G(\alpha) = F^{-1}(\alpha)$.

Quantile functions are generally set-valued when quantiles are not unique, although in the common case where the pdf does not contain any regions of 0 probability, the quantile function is the inverse of the cdf.

Example 1.19. The cdf of an exponential random variable is

$$F(x; \lambda) = 1 - \exp\left(-\frac{x}{\lambda}\right)$$

for $x \geq 0$ and $\lambda > 0$. Since $f(x; \lambda) > 0$ for $x > 0$, the quantile function is

$$F^{-1}(\alpha; \lambda) = -\lambda \ln(1 - \alpha).$$

The quantile function plays an important role in simulation of random variables. In particular, if $u \sim U(0, 1)$ ⁵, then $x = F^{-1}(u)$ is distributed F . For example, when u is a standard uniform ($U(0, 1)$), and $F^{-1}(\alpha)$ is the quantile function of an exponential random variable with shape parameter λ , then $x = F^{-1}(u; \lambda)$ follows an exponential(λ) distribution.

Theorem 1.5 (Probability Integral Transform). Let U be a standard uniform random variable, $F_X(x)$ be a continuous, increasing cdf. Then $\Pr(F^{-1}(U) < x) = F_X(x)$ and so $F^{-1}(U)$ is distributed F .

Proof. Let U be a standard uniform random variable, and for an $x \in R(X)$,

$$\Pr(U \leq F(x)) = F(x),$$

which follows from the definition of a standard uniform.

$$\begin{aligned} \Pr(U \leq F(x)) &= \Pr(F^{-1}(U) \leq F^{-1}(F(x))) \\ &= \Pr(F^{-1}(U) \leq x) \\ &= \Pr(X \leq x). \end{aligned}$$

□

⁵The mathematical notation \sim is read “distributed as”. For example, $x \sim U(0, 1)$ indicates that x is distributed as a standard uniform random variable.

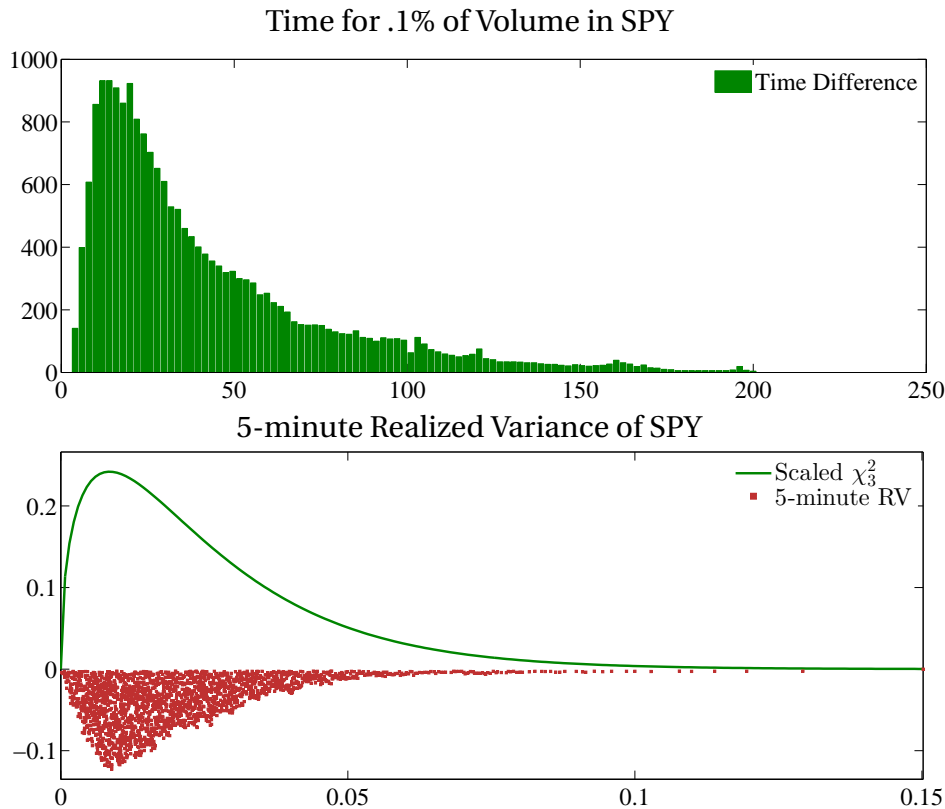


Figure 1.4: The left panel shows a histogram of the elapsed time in seconds required for .1% of the daily volume being traded to occur for SPY on May 31, 2012. The right panel shows both the fitted scaled χ^2 distribution and the raw data (mirrored below) for 5-minute “realized variance” estimates for SPY on May 31, 2012.

The key identity is that $\Pr(F^{-1}(U) \leq x) = \Pr(X \leq x)$, which shows that the distribution of $F^{-1}(U)$ is F by definition of the cdf. The right panel of figure 1.8 shows the relationship between the cdf of a standard normal and the associated quantile function. Applying $F(X)$ produces a uniform U through the cdf and applying $F^{-1}(U)$ produces X through the quantile function.

1.2.3 Common Univariate Distributions

Discrete

1.2.3.1 Bernoulli

A Bernoulli random variable is a discrete random variable which takes one of two values, 0 or 1. It is often used to model success or failure, where success is loosely defined. For example, a success may be the event that a trade was profitable net of costs, or the event that stock market volatility as measured by VIX was greater than 40%. The Bernoulli distribution depends on a single parameter p which determines the probability of success.

Parameters

$$p \in [0, 1]$$

Support

$$x \in \{0, 1\}$$

Probability Mass Function

$$f(x; p) = p^x (1 - p)^{1-x}, \quad p \geq 0$$

Moments

Mean	p
Variance	$p(1 - p)$

1.2.3.2 Poisson

A Poisson random variable is a discrete random variable taking values in $\{0, 1, \dots\}$. The Poisson depends on a single parameter λ (known as the intensity). Poisson random variables are often used to model counts of events during some interval, for example the number of trades executed over a 5-minute window.

Parameters

$$\lambda \geq 0$$

Support

$$x \in \{0, 1, \dots\}$$

Probability Mass Function

$$f(x; \lambda) = \frac{\lambda^x}{x!} \exp(-\lambda)$$

Moments

Mean	λ
Variance	λ

Continuous**1.2.3.3 Normal (Gaussian)**

The normal is the most important univariate distribution in financial economics. It is the familiar “bell-shaped” distribution, and is used heavily in hypothesis testing and in modeling (net) asset returns (e.g. $r_t = \ln P_t - \ln P_{t-1}$ or $r_t = \frac{P_t - P_{t-1}}{P_{t-1}}$ where P_t is the price of the asset in period t).

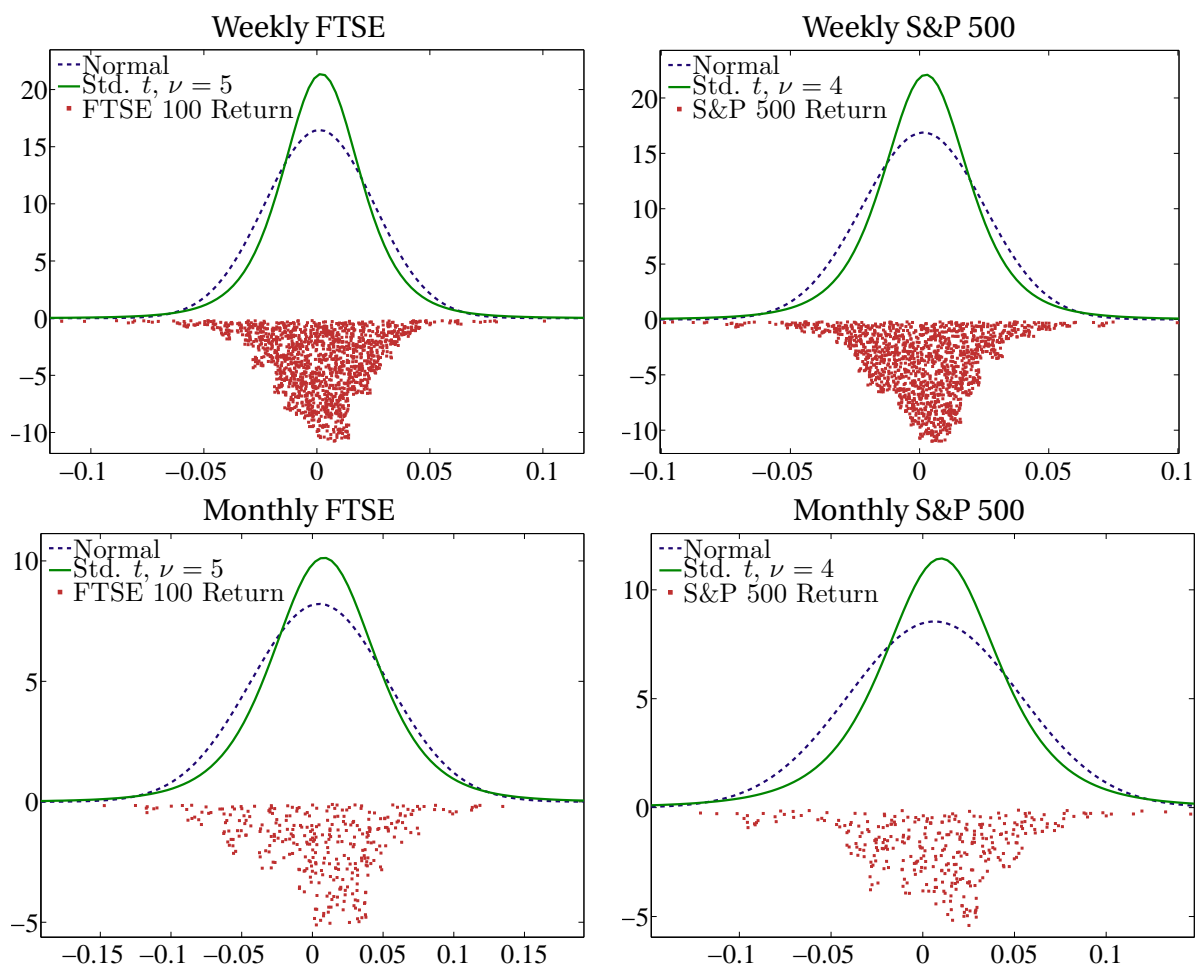


Figure 1.5: Weekly and monthly densities for the FTSE 100 and S&P 500. All panels plot the pdf of a normal and a standardized Student's t using parameters estimated with maximum likelihood estimation (See Chapter 1). The points below 0 on the y-axis show the actual returns observed during this period.

Parameters

$$\mu \in (-\infty, \infty), \sigma^2 \geq 0$$

Support

$$x \in (-\infty, \infty)$$

Probability Density Function

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

Cumulative Distribution Function

$$F(x; \mu, \sigma^2) = \frac{1}{2} + \frac{1}{2} \operatorname{erf}\left(\frac{1}{\sqrt{2}} \frac{x-\mu}{\sigma}\right) \text{ where erf is the error function.}^6$$

Moments

Mean	μ
Variance	σ^2
Median	μ
Skewness	0
Kurtosis	3

Notes

The normal with mean μ and variance σ^2 is written $N(\mu, \sigma^2)$. A normally distributed random variable with $\mu = 0$ and $\sigma^2 = 1$ is known as a standard normal. Figure 1.5 shows the fit normal distribution to the FTSE 100 and S&P 500 using both weekly and monthly returns for the period 1984–2012. Below each figure is a plot of the raw data.

1.2.3.4 Log-Normal

Log-normal random variables are closely related to normals. If X is log-normal, then $Y = \ln(X)$ is normal. Like the normal, the log-normal family depends on two parameters, μ and σ^2 , although unlike the normal these parameters do not correspond to the mean and variance. Log-normal random variables are commonly used to model gross returns, P_{t+1}/P_t (although it is often simpler to model $r_t = \ln P_t - \ln P_{t-1} = \ln(P_t/P_{t-1})$ which is normally distributed).

Parameters

$$\mu \in (-\infty, \infty), \sigma^2 \geq 0$$

⁶The error function does not have a closed form and is defined

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-s^2) \, ds.$$

Support

$$x \in (0, \infty)$$

Probability Density Function

$$f(x; \mu, \sigma^2) = \frac{1}{x\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\ln x - \mu)^2}{2\sigma^2}\right)$$

Cumulative Distribution Function

Since $Y = \ln(X) \sim N(\mu, \sigma^2)$, the cdf is the same as the normal only using $\ln x$ in place of x .

Moments

$$\begin{aligned} \text{Mean} & \exp\left(\mu + \frac{\sigma^2}{2}\right) \\ \text{Median} & \exp(\mu) \\ \text{Variance} & \{\exp(\sigma^2) - 1\} \exp(2\mu + \sigma^2) \end{aligned}$$

1.2.3.5 χ^2 (Chi-square)

χ^2_ν random variables depend on a single parameter ν known as the degree-of-freedom. They are commonly encountered when testing hypotheses, although they are also used to model continuous variables which are non-negative such as conditional variances. χ^2_ν random variables are closely related to standard normal random variables and are defined as the sum of ν independent standard normal random variables which have been squared. Suppose Z_1, \dots, Z_ν are standard normally distributed and independent, then $x = \sum_{i=1}^\nu z_i^2$ follows a χ^2_ν .⁷

Parameters

$$\nu \in [0, \infty)$$

Support

$$x \in [0, \infty)$$

Probability Density Function

$$f(x; \nu) = \frac{1}{2^{\frac{\nu}{2}} \Gamma(\frac{\nu}{2})} x^{\frac{\nu-2}{2}} \exp\left(-\frac{x}{2}\right), \quad \nu \in \{1, 2, \dots\} \text{ where } \Gamma(a) \text{ is the Gamma function.}^8$$

Cumulative Distribution Function

$$F(x; \nu) = \frac{1}{\Gamma(\frac{\nu}{2})} \gamma\left(\frac{\nu}{2}, \frac{x}{2}\right) \text{ where } \gamma(a, b) \text{ is the lower incomplete gamma function.}$$

⁷ ν does not need to be an integer,

⁸ The χ^2_ν is related to the gamma distribution which has pdf $f(x; \alpha, \beta) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} \exp(-x/\beta)$ by setting $\alpha = \nu/2$ and $\beta = 2$.

Moments

Mean	ν
Variance	2ν

Notes

Figure 1.4 shows a χ^2 pdf which was used to fit some simple estimators of the 5-minute variance of the S&P 500 from May 31, 2012. These were computed by summing and squaring 1-minute returns within a 5-minute interval (all using log prices). 5-minute variance estimators are important in high-frequency trading and other (slower) algorithmic trading.

1.2.3.6 Student's t and standardized Student's t

Student's t random variables are also commonly encountered in hypothesis testing and, like χ^2_ν random variables, are closely related to standard normals. Student's t random variables depend on a single parameter, ν , and can be constructed from two other independent random variables. If Z a standard normal, W a χ^2_ν and $Z \perp\!\!\!\perp W$, then $x = z/\sqrt{\frac{w}{\nu}}$ follows a Student's t distribution. Student's t are similar to normals except that they are heavier tailed, although as $\nu \rightarrow \infty$ a Student's t converges to a standard normal.

Support

$$x \in (-\infty, \infty)$$

Probability Density Function

$$f(x; \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}} \text{ where } \Gamma(a) \text{ is the Gamma function.}$$

Moments

Mean	$0, \nu > 1$
Median	0
Variance	$\frac{\nu}{\nu-2}, \nu > 2$
Skewness	$0, \nu > 3$
Kurtosis	$3\frac{(\nu-2)}{\nu-4}, \nu > 4$

Notes

When $\nu = 1$, a Student's t is known as a Cauchy random variable. Cauchy random variables are so heavy-tailed that even the mean does not exist.

The standardized Student's t extends the usual Student's t in two directions. First, it removes the variance's dependence on ν so that the scale of the random variable can be established separately from the degree of freedom parameter. Second, it explicitly adds location and scale parameters so that if Y is a Student's t random variable with degree of freedom ν , then

$$x = \mu + \sigma \frac{\sqrt{\nu-2}}{\sqrt{\nu}} y$$

follows a standardized Student's t distribution ($\nu > 2$ is required). The standardized Student's t is commonly used to model heavy-tailed return distributions such as stock market indices.

Figure 1.5 shows the fit (using maximum likelihood) standardized t distribution to the FTSE 100 and S&P 500 using both weekly and monthly returns from the period 1984–2012. The typical degree of freedom parameter was around 4, indicating that (unconditional) distributions are heavy-tailed with a large kurtosis.

1.2.3.7 Uniform

The continuous uniform is commonly encountered in certain test statistics, especially those testing whether assumed densities are appropriate for a particular series. Uniform random variables, when combined with quantile functions, are also useful for simulating random variables.

Parameters

a, b the end points of the interval, where $a < b$

Support

$$x \in [a, b]$$

Probability Density Function

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

Cumulative Distribution Function

$$F(x) = \frac{x-a}{b-a} \text{ for } a \leq x \leq b, F(x) = 0 \text{ for } x < a \text{ and } F(x) = 1 \text{ for } x > b$$

Moments

Mean	$\frac{b+a}{2}$
Median	$\frac{b+a}{2}$
Variance	$\frac{(b-a)^2}{12}$
Skewness	0
Kurtosis	$\frac{9}{5}$

Notes

A standard uniform has $a = 0$ and $b = 1$. When $x \sim F$, then $F(x) \sim U(0, 1)$

1.3 Multivariate Random Variables

While univariate random variables are very important in financial economics, most applications require the use multivariate random variables. Multivariate random variables allow the relation-

ship between two or more random quantities to be modeled and studied. For example, the joint distribution of equity and bond returns is important for many investors.

Throughout this section, the multivariate random variable is assumed to have n components,

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$$

which are arranged into a column vector. The definition of a multivariate random variable is virtually identical to that of a univariate random variable, only mapping $\omega \in \Omega$ to the n -dimensional space \mathbb{R}^n .

Definition 1.23 (Multivariate Random Variable). Let (Ω, \mathcal{F}, P) be a probability space. If $X : \Omega \rightarrow \mathbb{R}^n$ is a real-valued vector function having its domain the elements of Ω , then $X : \Omega \rightarrow \mathbb{R}^n$ is called a (multivariate) n -dimensional random variable.

Multivariate random variables, like univariate random variables, are technically functions of events in the underlying probability space $X(\omega)$, although the function argument ω (the event) is usually suppressed.

Multivariate random variables can be either discrete or continuous. Discrete multivariate random variables are fairly uncommon in financial economics and so the remainder of the chapter focuses exclusively on the continuous case. The characterization of a what makes a multivariate random variable continuous is also virtually identical to that in the univariate case.

Definition 1.24 (Continuous Multivariate Random Variable). A multivariate random variable is said to be continuous if its range is uncountably infinite and if there exists a non-negative valued function $f(x_1, \dots, x_n)$ defined for all $(x_1, \dots, x_n) \in \mathbb{R}^n$ such that for any event $B \subset R(X)$,

$$\Pr(B) = \int \dots \int_{\{x_1, \dots, x_n\} \in B} f(x_1, \dots, x_n) dx_1 \dots dx_n \quad (1.12)$$

and $f(x_1, \dots, x_n) = 0$ for all $(x_1, \dots, x_n) \notin R(X)$.

Multivariate random variables, at least when continuous, are often described by their probability density function.

Definition 1.25 (Continuous Density Function Characterization). A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a member of the class of multivariate continuous density functions if and only if $f(x_1, \dots, x_n) \geq 0$ for all $x \in \mathbb{R}^n$ and

$$\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_1 \dots dx_n = 1. \quad (1.13)$$

Definition 1.26 (Multivariate Probability Density Function). The function $f(x_1, \dots, x_n)$ is called a multivariate probability function (pdf).

A multivariate density, like a univariate density, is a function which is everywhere non-negative and which integrates to unity. Figure 1.7 shows the fit joint probability density function to weekly returns on the FTSE 100 and S&P 500 (assuming that returns are normally distributed). Two views are presented – one shows the 3-dimensional plot of the pdf and the other shows the iso-probability contours of the pdf. The figure also contains a scatter plot of the raw weekly data for comparison. All parameters were estimated using maximum likelihood.

Example 1.20. Suppose X is a bivariate random variable, then the function $f(x_1, x_2) = \frac{3}{2}(x_1^2 + x_2^2)$ defined on $[0, 1] \times [0, 1]$ is a valid probability density function.

Example 1.21. Suppose X is a bivariate standard normal random variable. Then the probability density function of X is

$$f(x_1, x_2) = \frac{1}{2\pi} \exp\left(-\frac{x_1^2 + x_2^2}{2}\right).$$

The multivariate cumulative distribution function is virtually identical to that in the univariate case, and measure the total probability between $-\infty$ (for each element of X) and some point.

Definition 1.27 (Multivariate Cumulative Distribution Function). The joint cumulative distribution function of an n -dimensional random variable X is defined by

$$F(x_1, \dots, x_n) = \Pr(X_i \leq x_i, i = 1, \dots, n)$$

for all $(x_1, \dots, x_n) \in \mathbb{R}^n$, and is given by

$$F(x_1, \dots, x_n) = \int_{-\infty}^{x_n} \dots \int_{-\infty}^{x_1} f(s_1, \dots, s_n) ds_1 \dots ds_n. \quad (1.14)$$

Example 1.22. Suppose X is a bivariate random variable with probability density function

$$f(x_1, x_2) = \frac{3}{2}(x_1^2 + x_2^2)$$

defined on $[0, 1] \times [0, 1]$. Then the associated cdf is

$$F(x_1, x_2) = \frac{x_1^3 x_2 + x_1 x_2^3}{2}.$$

Figure 1.6 shows the joint cdf of the density in the previous example. As was the case for univariate random variables, the probability density function can be determined by differentiating the cumulative distribution function with respect to each component.

Theorem 1.6 (Relationship between cdf and pdf). Let $f(x_1, \dots, x_n)$ and $F(x_1, \dots, x_n)$ represent the pdf and cdf of an n -dimensional continuous random variable X , respectively. The density function for X can be defined as $f(x_1, \dots, x_n) = \frac{\partial^n F(\mathbf{x})}{\partial x_1 \partial x_2 \dots \partial x_n}$ whenever $f(x_1, \dots, x_n)$ is continuous and $f(x_1, \dots, x_n) = 0$ elsewhere.

Example 1.23. Suppose X is a bivariate random variable with cumulative distribution function $F(x_1, x_2) = \frac{x_1^3 x_2 + x_1 x_2^3}{2}$. The probability density function can be determined using

$$\begin{aligned} f(x_1, x_2) &= \frac{\partial^2 F(x_1, x_2)}{\partial x_1 \partial x_2} \\ &= \frac{1}{2} \frac{\partial (3x_1^2 x_2 + x_2^3)}{\partial x_2} \\ &= \frac{3}{2} (x_1^2 + x_2^2). \end{aligned}$$

1.3.1 Marginal Densities and Distributions

The marginal distribution is the first concept unique to multivariate random variables. Marginal densities and distribution functions summarize the information in a subset, usually a single component, of X by averaging over all possible values of the components of X which are not being marginalized. This involves integrating out the variables which are not of interest. First, consider the bivariate case.

Definition 1.28 (Bivariate Marginal Probability Density Function). Let X be a bivariate random variable comprised of X_1 and X_2 . The marginal distribution of X_1 is given by

$$f_1(x_1) = \int_{-\infty}^{\infty} f(x_1, x_2) dx_2. \quad (1.15)$$

The marginal density of X_1 is a density function where X_2 has been integrated out. This integration is simply a form of averaging – varying x_2 according to the probability associated with each value of x_2 – and so the marginal is only a function of x_1 . Both probability density functions and cumulative distribution functions have marginal versions.

Example 1.24. Suppose X is a bivariate random variable with probability density function

$$f(x_1, x_2) = \frac{3}{2} (x_1^2 + x_2^2)$$

and is defined on $[0, 1] \times [0, 1]$. The marginal probability density function for X_1 is

$$f_1(x_1) = \frac{3}{2} \left(x_1^2 + \frac{1}{3} \right),$$

and by symmetry the marginal probability density function of X_2 is

$$f_2(x_2) = \frac{3}{2} \left(x_2^2 + \frac{1}{3} \right).$$

Example 1.25. Suppose X is a bivariate random variable with probability density function $f(x_1, x_2) = 6(x_1 x_2^2)$ and is defined on $[0, 1] \times [0, 1]$. The marginal probability density functions for X_1 and X_2 are

$$f_1(x_1) = 2x_1 \text{ and } f_2(x_2) = 3x_2^2.$$

Example 1.26. Suppose X is bivariate normal with parameters $\mu = [\mu_1 \mu_2]'$ and

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix},$$

then the marginal pdf of X_1 is $N(\mu_1, \sigma_1^2)$, and the marginal pdf of X_2 is $N(\mu_2, \sigma_2^2)$.

Figure 1.7 shows the fit marginal distributions to weekly returns on the FTSE 100 and S&P 500 assuming that returns are normally distributed. Marginal pdfs can be transformed into marginal cdfs through integration.

Definition 1.29 (Bivariate Marginal Cumulative Distribution Function). The cumulative marginal distribution function of X_1 in bivariate random variable X is defined by

$$F_1(x_1) = \Pr(X_1 \leq x_1)$$

for all $x_1 \in \mathbb{R}$, and is given by

$$F_1(x_1) = \int_{-\infty}^{x_1} f_1(s_1) ds_1.$$

The general j -dimensional marginal distribution partitions the n -dimensional random variable X into two blocks, and constructs the marginal distribution for the first j by integrating out (averaging over) the remaining $n - j$ components of X . In the definition, both X_1 and X_2 are vectors.

Definition 1.30 (Marginal Probability Density Function). Let X be a n -dimensional random variable and partition the first j ($1 \leq j < n$) elements of X into X_1 , and the remainder into X_2 so that $X = [X_1' X_2']'$. The marginal probability density function for X_1 is given by

$$f_{1,\dots,j}(x_1, \dots, x_j) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_{j+1} \dots dx_n. \quad (1.16)$$

The marginal cumulative distribution function is related to the marginal probability density function in the same manner as the joint probability density function is related to the cumulative distribution function. It also has the same interpretation.

Definition 1.31 (Marginal Cumulative Distribution Function). Let X be a n -dimensional random variable and partition the first j ($1 \leq j < n$) elements of X into X_1 , and the remainder into X_2 so that $X = [X_1' X_2']'$. The marginal cumulative distribution function for X_1 is given by

$$F_{1,\dots,j}(x_1, \dots, x_j) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_j} f_{1,\dots,j}(s_1, \dots, s_j) ds_1 \dots ds_j. \quad (1.17)$$

1.3.2 Conditional Distributions

Marginal distributions provide the tools needed to model the distribution of a subset of the components of a random variable while averaging over the other components. Conditional densities and distributions, on the other hand, consider a subset of the components a random variable conditional on observing a specific value for the remaining components. In practice, the vast majority of modeling makes use of conditioning information where the interest is in understanding the distribution of a random variable conditional on the observed values of some other random variables. For example, consider the problem of modeling the expected return of an individual stock. Balance sheet information such as the book value of assets, earnings and return on equity are all available, and can be conditioned on to model the conditional distribution of the stock's return.

First, consider the bivariate case.

Definition 1.32 (Bivariate Conditional Probability Density Function). Let X be a bivariate random variable comprised of X_1 and X_2 . The conditional probability density function for X_1 given that $X_2 \in B$ where B is an event where $\Pr(X_2 \in B) > 0$ is

$$f(x_1 | X_2 \in B) = \frac{\int_B f(x_1, x_2) dx_2}{\int_B f_2(x_2) dx_2}. \quad (1.18)$$

When B is an elementary event (e.g. single point), so that $\Pr(X_2 = x_2) = 0$ and $f_2(x_2) > 0$, then

$$f(x_1 | X_2 = x_2) = \frac{f(x_1, x_2)}{f_2(x_2)}. \quad (1.19)$$

Conditional density functions differ slightly depending on whether the conditioning variable is restricted to a set or a point. When the conditioning variable is specified to be a set where $\Pr(X_2 \in B) > 0$, then the conditional density is the joint probability of X_1 and $X_2 \in B$ divided by the marginal probability of $X_2 \in B$. When the conditioning variable is restricted to a point, the conditional density is the ratio of the joint pdf to the margin pdf of X_2 .

Example 1.27. Suppose X is a bivariate random variable with probability density function

$$f(x_1, x_2) = \frac{3}{2} (x_1^2 + x_2^2)$$

and is defined on $[0, 1] \times [0, 1]$. The conditional probability of X_1 given $X_2 \in [\frac{1}{2}, 1]$

$$f\left(x_1 | X_2 \in \left[\frac{1}{2}, 1\right]\right) = \frac{1}{11} (12x_1^2 + 7),$$

the conditional probability density function of X_1 given $X_2 \in [0, \frac{1}{2}]$ is

$$f\left(x_1 | X_2 \in \left[0, \frac{1}{2}\right]\right) = \frac{1}{5} (12x_1^2 + 1),$$

and the conditional probability density function of X_1 given $X_2 = x_2$ is

$$f(x_1 | X_2 = x_2) = \frac{x_1^2 + x_2^2}{x_2^2 + 1}.$$

Figure 1.6 shows the joint pdf along with both types of conditional densities. The upper left panel shows that conditional density for $X_2 \in [0.25, 0.5]$. The highlighted region contains the components of the joint pdf which are averaged to produce the conditional density. The lower left also shows the pdf but also shows three (non-normalized) conditional densities of the form $f(x_1 | x_2)$. The lower right pane shows these three densities correctly normalized.

The previous example shows that, in general, the conditional probability density function differs as the region used changes.

Example 1.28. Suppose X is bivariate normal with mean $\mu = [\mu_1 \mu_2]'$ and covariance

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix},$$

then the conditional distribution of X_1 given $X_2 = x_2$ is $N\left(\mu_1 + \frac{\sigma_{12}}{\sigma_2^2} (x_2 - \mu_2), \sigma_1^2 - \frac{\sigma_{12}^2}{\sigma_2^2}\right)$.

Marginal distributions and conditional distributions are related in a number of ways. One obvious way is that $f(x_1|X_2 \in R(X_2)) = f_1(x_1)$ – that is, the conditional probability of X_1 given that X_2 is in its range is the marginal pdf of X_1 . This holds since integrating over all values of x_2 is essentially not conditioning on anything (which is known as the unconditional, and a marginal density could, in principle, be called the unconditional density since it averages across all values of the other variable).

The general definition allows for an n -dimensional random vector where the conditioning variable has a dimension between 1 and $j < n$.

Definition 1.33 (Conditional Probability Density Function). Let $f(x_1, \dots, x_n)$ be the joint density function for an n -dimensional random variable $X = [X_1 \dots X_n]'$ and partition the first j ($1 \leq j < n$) elements of X into X_1 , and the remainder into X_2 so that $X = [X_1' X_2']'$. The conditional probability density function for X_1 given that $X_2 \in B$ is given by

$$f(x_1, \dots, x_j | X_2 \in B) = \frac{\int_{(x_{j+1}, \dots, x_n) \in B} f(x_1, \dots, x_n) dx_n \dots dx_{j+1}}{\int_{(x_{j+1}, \dots, x_n) \in B} f_{j+1, \dots, n}(x_{j+1}, \dots, x_n) dx_n \dots dx_{j+1}}, \quad (1.20)$$

and when B is an elementary event (denoted \mathbf{x}_2) and if $f_{j+1, \dots, n}(\mathbf{x}_2) > 0$,

$$f(x_1, \dots, x_j | X_2 = \mathbf{x}_2) = \frac{f(x_1, \dots, x_j, \mathbf{x}_2)}{f_{j+1, \dots, n}(\mathbf{x}_2)} \quad (1.21)$$

In general the simplified notation $f(x_1, \dots, x_j | \mathbf{x}_2)$ will be used to represent $f(x_1, \dots, x_j | X_2 = \mathbf{x}_2)$.

1.3.3 Independence

A special relationship exists between the joint probability density function and the marginal density functions when random variables are independent– the joint must be the product of each marginal.

Theorem 1.7 (Independence of Random Variables). *The random variables X_1, \dots, X_n with joint density function $f(x_1, \dots, x_n)$ are independent if and only if*

$$f(x_1, \dots, x_n) = \prod_{i=1}^n f_i(x_i) \quad (1.22)$$

where $f_i(x_i)$ is the marginal distribution of X_i .

The intuition behind this result follows from the fact that when the components of a random variable are independent, any change in one component has no information for the others. In other words, both marginals and conditionals must be the same.

Example 1.29. Let X be a bivariate random variable with probability density function $f(x_1, x_2) = x_1 x_2$ on $[0, 1] \times [0, 1]$, then X_1 and X_2 are independent. This can be verified since

$$f_1(x_1) = x_1 \text{ and } f_2(x_2) = x_2$$

so that the joint is the product of the two marginal densities.

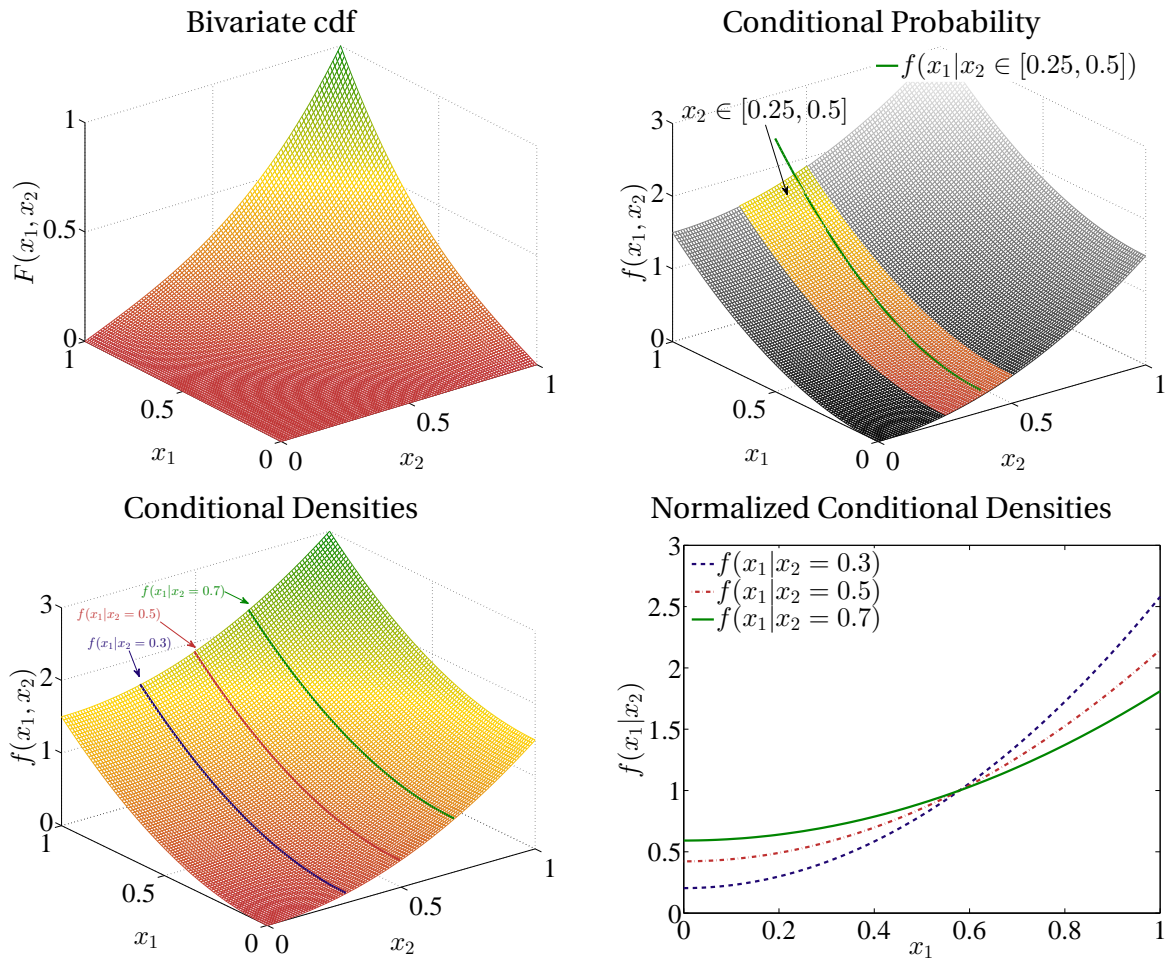


Figure 1.6: These four panels show four views of a distribution defined on $[0, 1] \times [0, 1]$. The upper left panel shows the joint cdf. The upper right shows the pdf along with the portion of the pdf used to construct a conditional distribution $f(x_1 | x_2 \in [0.25, 0.5])$. The line shows the actual correctly scaled conditional distribution which is only a function of x_1 plotted at $E[X_2 | X_2 \in [0.25, 0.5]]$. The lower left panel also shows the pdf along with three non-normalized conditional densities. The bottom right panel shows the correctly normalized conditional densities.

Independence is a very strong concept, and it carries over from random variables to functions of random variables as long as each function involves only one random variable.⁹

Theorem 1.8 (Independence of Functions of Independent Random Variables). *Let X_1 and X_2 be independent random variables and define $y_1 = Y_1(x_1)$ and $y_2 = Y_2(x_2)$, then the random variables Y_1 and Y_2 are independent.*

Independence is often combined with an assumption that the marginal distribution is the same to simplify the analysis of collections of random data.

Definition 1.34 (Independent, Identically Distributed). Let $\{X_i\}$ be a sequence of random variables. If the marginal distribution for X_i is the same for all i and $X_i \perp\!\!\!\perp X_j$ for all $i \neq j$, then $\{X_i\}$ is said to be an independent, identically distributed (i.i.d.) sequence.

1.3.4 Bayes Rule

Bayes rule is used both in financial economics and econometrics. In financial economics, it is often used to model agents learning, and in econometrics it is used to make inference about unknown parameters given observed data (a branch known as Bayesian econometrics). Bayes rule follows directly from the definition of a conditional density so that the joint can be factored into a conditional and a marginal. Suppose X is a bivariate random variable, then

$$\begin{aligned} f(x_1, x_2) &= f(x_1|x_2) f_2(x_2) \\ &= f(x_2|x_1) f_1(x_1). \end{aligned}$$

The joint can be factored two ways, and equating the two factorizations results in Bayes rule.

Definition 1.35 (Bivariate Bayes Rule). Let X be a bivariate random variable with components X_1 and X_2 , then

$$f(x_1|x_2) = \frac{f(x_2|x_1) f_1(x_1)}{f_2(x_2)} \quad (1.23)$$

Bayes rule states that the probability of observing X_1 given a value of X_2 is equal to the joint probability of the two random variables divided by the marginal probability of observing X_2 . Bayes rule is normally applied where there is a belief about X_1 ($f_1(x_1)$, called a *prior*), and the conditional distribution of X_1 given X_2 is a known density ($f(x_2|x_1)$, called the *likelihood*), which combine to form a belief about X_1 ($f(x_1|x_2)$, called the *posterior*). The marginal density of X_2 is not important when using Bayes rule since the numerator is still proportional to the conditional density of X_1 given X_2 since $f_2(x_2)$ is a number, and so it is common to express the non-normalized posterior as

$$f(x_1|x_2) \propto f(x_2|x_1) f_1(x_1),$$

where \propto is read “is proportional to”.

⁹This can be generalized to the full multivariate case where X is an n -dimensional random variable where the first j components are independent from the last $n - j$ components defining $y_1 = Y_1(x_1, \dots, x_j)$ and $y_2 = Y_2(x_{j+1}, \dots, x_n)$.

Example 1.30. Suppose interest lies in the probability a firm does bankrupt which can be modeled as a Bernoulli distribution. The parameter p is unknown but, given a value of p , the likelihood that a firm goes bankrupt is

$$f(x|p) = p^x (1-p)^{1-x}.$$

While p is known, a prior for the bankruptcy rate can be specified. Suppose the prior for p follows a Beta(α, β) distribution which has pdf

$$f(p) = \frac{p^{\alpha-1} (1-p)^{\beta-1}}{B(\alpha, \beta)}$$

where $B(a, b)$ is Beta function that acts as a normalizing constant.¹⁰ The Beta distribution has support on $[0, 1]$ and nests the standard uniform as a special case when $\alpha = \beta = 1$. The expected value of a random variable with a Beta(α, β) is $\frac{\alpha}{\alpha+\beta}$ and the variance is $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ where $\alpha > 0$ and $\beta > 0$.

Using Bayes rule,

$$\begin{aligned} f(p|x) &\propto p^x (1-p)^{1-x} \times \frac{p^{\alpha-1} (1-p)^{\beta-1}}{B(\alpha, \beta)} \\ &= \frac{p^{\alpha-1+x} (1-p)^{\beta-x}}{B(\alpha, \beta)}. \end{aligned}$$

Note that this isn't a density since it has the wrong normalizing constant. However, the component of the density which contains p is $p^{(\alpha+x)-1} (1-p)^{(\beta-x+1)-1}$ (known as the *kernel*) is the same as in the Beta distribution, only with different parameters. Thus the posterior, $f(p|x)$ is Beta($\alpha+x, \beta-x+1$). Since the posterior is the same as the prior, it could be combined with another observation (and the Bernoulli likelihood) to produce an updated posterior. When a Bayesian problem has this property, the prior density said to be conjugate to the likelihood.

Example 1.31. Suppose M is a random variable representing the score on the midterm, and interest lies in the final course grade, C . The prior for C is normal with mean μ and variance σ^2 , and that the distribution of M given C is also conditionally normal with mean C and variance τ^2 . Bayes rule can be used to make inference on the final course grade given the midterm grade.

$$\begin{aligned} f(c|m) &\propto f(m|c) f_C(c) \\ &\propto \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{(m-c)^2}{2\tau^2}\right) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(c-\mu)^2}{2\sigma^2}\right) \\ &= K \exp\left(-\frac{1}{2} \left\{ \frac{(m-c)^2}{\tau^2} + \frac{(c-\mu)^2}{\sigma^2} \right\}\right) \\ &= K \exp\left(-\frac{1}{2} \left\{ \frac{c^2}{\tau^2} + \frac{c^2}{\sigma^2} - \frac{2cm}{\tau^2} - \frac{2c\mu}{\sigma^2} + \frac{m^2}{\tau^2} + \frac{\mu^2}{\sigma^2} \right\}\right) \\ &= K \exp\left(-\frac{1}{2} \left\{ c^2 \left(\frac{1}{\tau^2} + \frac{1}{\sigma^2} \right) - 2c \left(\frac{m}{\tau^2} + \frac{\mu}{\sigma^2} \right) + \left(\frac{m^2}{\tau^2} + \frac{\mu^2}{\sigma^2} \right) \right\}\right) \end{aligned}$$

¹⁰The beta function can only be given as an indefinite integral,

$$B(a, b) = \int_0^1 s^{a-1} (1-s)^{b-1} ds.$$

This (non-normalized) density can be shown to have the kernel of a normal by completing the square,¹¹

$$f(c|m) \propto \exp \left(-\frac{1}{2 \left(\frac{1}{\tau^2} + \frac{1}{\sigma^2} \right)^{-1}} \left(c - \frac{\left(\frac{m}{\tau^2} + \frac{\mu}{\sigma^2} \right)}{\left(\frac{1}{\tau^2} + \frac{1}{\sigma^2} \right)} \right)^2 \right).$$

This is the kernel of a normal density with mean

$$\frac{\left(\frac{m}{\tau^2} + \frac{\mu}{\sigma^2} \right)}{\left(\frac{1}{\tau^2} + \frac{1}{\sigma^2} \right)},$$

and variance

$$\left(\frac{1}{\tau^2} + \frac{1}{\sigma^2} \right)^{-1}.$$

The mean is a weighted average of the prior mean, μ and the midterm score, m , where the weights are determined by the inverse variance of the prior and conditional distributions. Since the weights are proportional to the inverse of the variance, a small variance leads to a relatively large weight. If $\tau^2 = \sigma^2$, then the posterior mean is the average of the prior mean and the midterm score. The variance of the posterior depends on the uncertainty in the prior (σ^2) and the uncertainty in the data (τ^2). The posterior variance is always less than the smaller of σ^2 and τ^2 . Like the Bernoulli-Beta combination in the previous problem, the normal distribution is a conjugate prior when the conditional density is normal.

1.3.5 Common Multivariate Distributions

1.3.5.1 Multivariate Normal

Like the univariate normal, the multivariate normal depends on 2 parameters, μ and n by 1 vector of means and Σ an n by n positive semi-definite covariance matrix. The multivariate normal is closed to both to marginalization and conditioning – in other words, if X is multivariate normal, then all marginal distributions of X are normal, and so are all conditional distributions of X_1 given X_2 for any partitioning.

Parameters

$\mu \in \mathbb{R}^n$, Σ a positive semi-definite matrix

Support

$\mathbf{x} \in \mathbb{R}^n$

¹¹ Suppose a quadratic in x has the form $ax^2 + bx + c$. Then

$$ax^2 + bx + c = a(x - d)^2 + e$$

where $d = b/(2a)$ and $e = c - b^2/(4a)$.

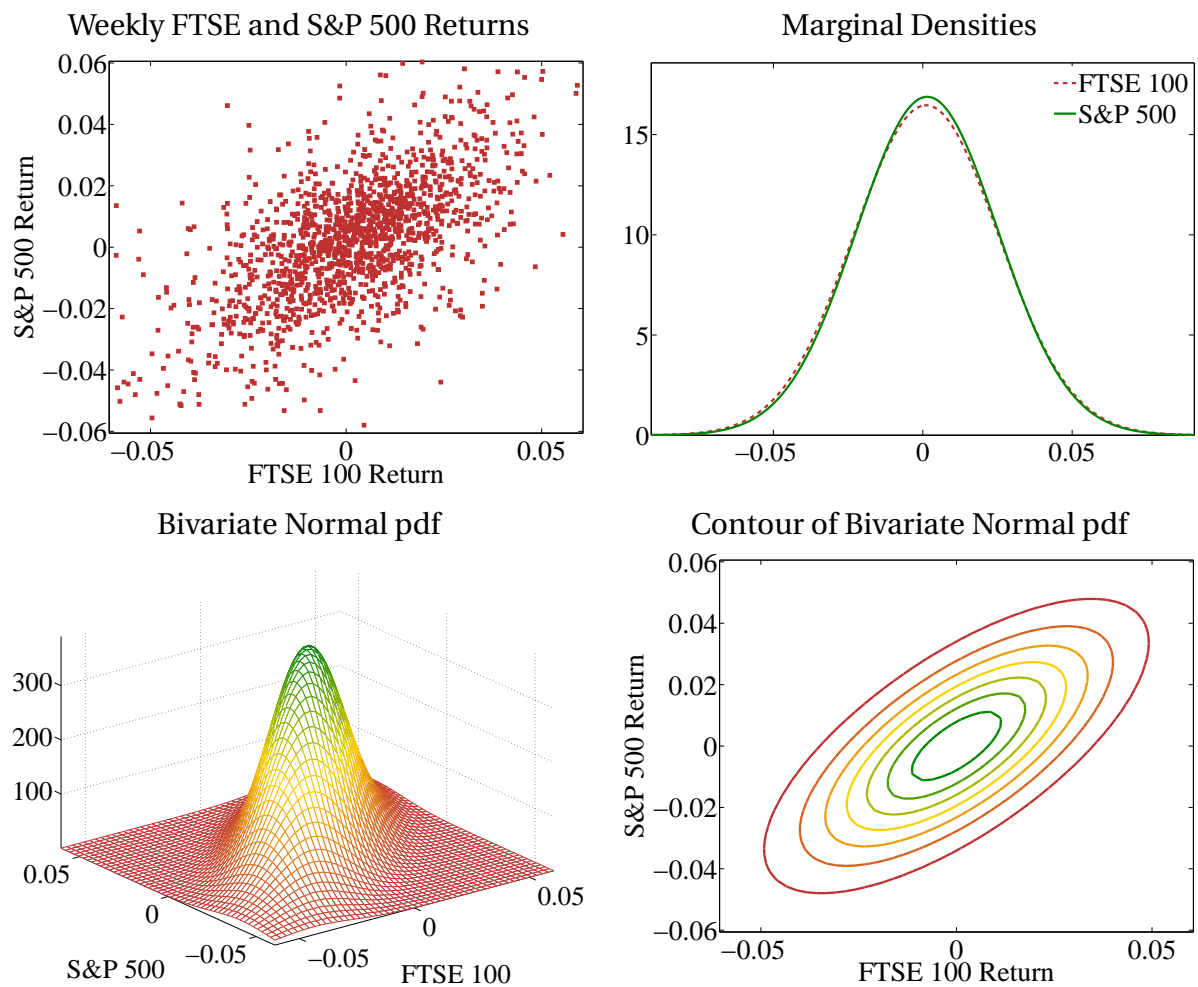


Figure 1.7: These four figures show different views of the weekly returns of the FTSE 100 and the S&P 500. The top left contains a scatter plot of the raw data. The top right shows the marginal distributions from a fit bivariate normal distribution (using maximum likelihood). The bottom two panels show two representations of the joint probability density function.

Probability Density Function

$$f(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

Cumulative Distribution Function

Can be expressed as a series of n univariate normal cdfs using repeated conditioning.

Moments

Mean	$\boldsymbol{\mu}$
Median	$\boldsymbol{\mu}$
Variance	$\boldsymbol{\Sigma}$
Skewness	0
Kurtosis	3

Marginal Distribution

The marginal distribution for the first j components is

$$f_{X_1, \dots, X_j}(x_1, \dots, x_j) = (2\pi)^{-\frac{j}{2}} |\boldsymbol{\Sigma}_{11}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}_1 - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}_{11}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1)\right),$$

where it is assumed that the marginal distribution is that of the first j random variables¹², $\boldsymbol{\mu} = [\boldsymbol{\mu}_1' \boldsymbol{\mu}_2']'$ where $\boldsymbol{\mu}_1$ correspond to the first j entries, and

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}' & \boldsymbol{\Sigma}_{22} \end{bmatrix}.$$

In other words, the distribution of $[X_1, \dots, X_j]'$ is $N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$. Moreover, the marginal distribution of a single element of X is $N(\mu_i, \sigma_i^2)$ where μ_i is the i^{th} element of $\boldsymbol{\mu}$ and σ_i^2 is the i^{th} diagonal element of $\boldsymbol{\Sigma}$.

Conditional Distribution

The conditional probability of X_1 given $X_2 = \mathbf{x}_2$ is

$$N(\boldsymbol{\mu}_1 + \boldsymbol{\beta}'(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\beta}'\boldsymbol{\Sigma}_{22}\boldsymbol{\beta})$$

where $\boldsymbol{\beta} = \boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{12}'$.

When X is a bivariate normal random variable,

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}\right),$$

¹²Any two variables can be reordered in a multivariate normal by swapping their means and reordering the corresponding rows and columns of the covariance matrix.

the conditional distribution is

$$X_1|X_2 = x_2 \sim N\left(\mu_1 + \frac{\sigma_{12}}{\sigma_2^2}(x_2 - \mu_2), \sigma_1^2 - \frac{\sigma_{12}^2}{\sigma_2^2}\right),$$

where the variance can be seen to always be positive since $\sigma_1^2\sigma_2^2 \geq \sigma_{12}^2$ by the Cauchy-Schwarz inequality (see 1.15).

Notes

The multivariate Normal has a number of novel and useful properties:

- A standard multivariate normal has $\boldsymbol{\mu} = \mathbf{0}$ and $\boldsymbol{\Sigma} = \mathbf{I}_n$.
- If the covariance between elements i and j equals zero (so that $\sigma_{ij} = 0$), they are independent.
- For the normal, zero covariance (or correlation) implies independence. This is not true of most other multivariate random variables.
- Weighted sums of multivariate normal random variables are normal. In particular if \mathbf{c} is a n by 1 vector of weights, then $Y = \mathbf{c}'X$ is normal with mean $\mathbf{c}'\boldsymbol{\mu}$ and variance $\mathbf{c}'\boldsymbol{\Sigma}\mathbf{c}$.

1.4 Expectations and Moments

Expectations and moments are (non-random) functions of random variables that are useful in both understanding properties of random variables – e.g. when comparing the dispersion between two distributions – and when estimating parameters using a technique known as the method of moments (see Chapter 1).

1.4.1 Expectations

The expectation is the value, on average, of a random variable (or function of a random variable). Unlike common English language usage, where one's expectation is not well defined (e.g. could be the mean or the mode, another measure of the tendency of a random variable), the expectation in a probabilistic sense *always* averages over the possible values weighting by the probability of observing each value. The form of an expectation in the discrete case is particularly simple.

Definition 1.36 (Expectation of a Discrete Random Variable). The expectation of a discrete random variable, defined $E[X] = \sum_{x \in R(X)} x f(x)$, exists if and only if $\sum_{x \in R(X)} |x| f(x) < \infty$.

When the range of X is finite then the expectation always exists. When the range is infinite, such as when a random variable takes on values in the range $0, 1, 2, \dots$, the probability mass function must be sufficiently small for large values of the random variable in order for the expectation to exist.¹³ Expectations of continuous random variables are virtually identical, only replacing the sum with an integral.

¹³An expectation is said to be nonexistent when the sum converges to $\pm\infty$ or oscillates. The use of the $|x|$ in the definition of existence is to rule out both the $-\infty$ and the oscillating cases.

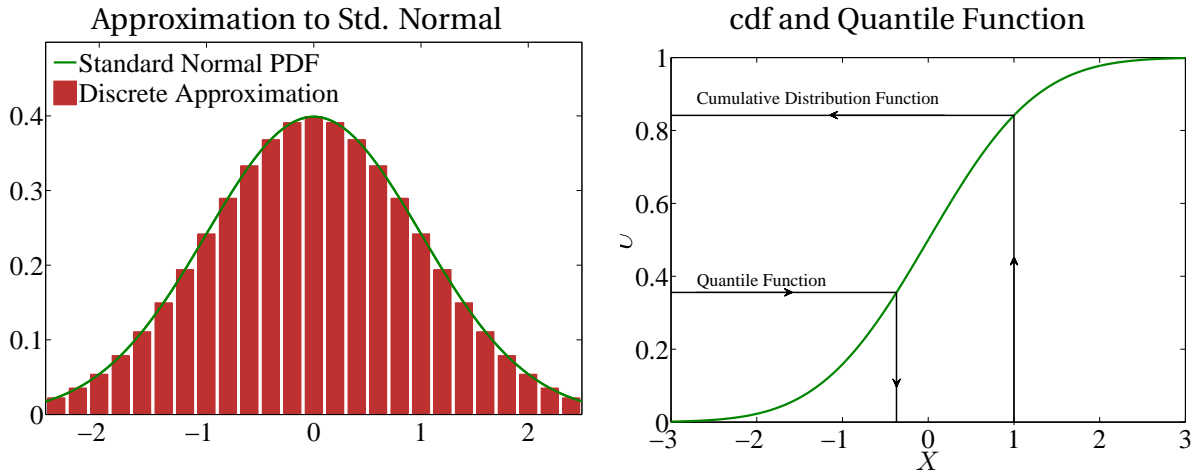


Figure 1.8: The left panel shows a standard normal and a discrete approximation. Discrete approximations are useful for approximating integrals in expectations. The right panel shows the relationship between the quantile function and the cdf.

Definition 1.37 (Expectation of a Continuous Random Variable). The expectation of a continuous random variable, defined $E[X] = \int_{-\infty}^{\infty} x f(x) dx$, exists if and only if $\int_{-\infty}^{\infty} |x| f(x) dx < \infty$.

The existence of an expectation is a somewhat difficult concept. For continuous random variables, expectations may not exist if the probability of observing an arbitrarily large value (in the absolute sense) is very high. For example, in a Student's t distribution when the degree of freedom parameter ν is 1 (also known as a Cauchy distribution), the probability of observing a value with size $|x|$ is proportional to x^{-1} for large x (in other words, $f(x) \propto c x^{-1}$) so that $x f(x) \approx c$ for large x . The range is unbounded, and so the integral of a constant, even if very small, will not converge, and so the expectation does not exist. On the other hand, when a random variable is bounded, its expectation always exists.

Theorem 1.9 (Expectation Existence for Bounded Random Variables). If $|x| < c$ for all $x \in R(X)$, then $E[X]$ exists.

The expectation operator, $E[\cdot]$ is generally defined for arbitrary functions of a random variable, $g(x)$. In practice, $g(x)$ takes many forms – x , x^2 , x^p for some p , $\exp(x)$ or something more complicated. Discrete and continuous expectations are closely related. Figure 1.8 shows a standard normal along with a discrete approximation where each bin has a width of 0.20 and the height is based on the pdf value at the mid-point of the bin. Treating the normal as a discrete distribution based on this approximation would provide reasonable approximations to the correct (integral) expectations.

Definition 1.38 (Expectation of a Function of Random Variable). The expectation of a random variable defined as a function of X , $Y = g(x)$, is $E[Y] = E[g(X)] = \int_{-\infty}^{\infty} g(x) f(x) dx$ exists if and only if $\int_{-\infty}^{\infty} |g(x)| dx < \infty$.

When $g(x)$ is either concave or convex, Jensen's inequality provides a relationship between the expected value of the function and the function of the expected value of the underlying random variable.

Theorem 1.10 (Jensen's Inequality). *If $g(\cdot)$ is a continuous convex function on an open interval containing the range of X , then $E[g(X)] \geq g(E[X])$. Similarly, if $g(\cdot)$ is a continuous concave function on an open interval containing the range of X , then $E[g(X)] \leq g(E[X])$.*

The inequalities become strict if the functions are strictly convex (or concave) as long as X is not degenerate.¹⁴ Jensen's inequality is common in economic applications. For example, standard utility functions ($U(\cdot)$) are assumed to be concave which reflects the idea that marginal utility ($U'(\cdot)$) is decreasing in consumption (or wealth). Applying Jensen's inequality shows that if consumption is random, then $E[U(c)] < U(E[c])$ – in other words, the economic agent is worse off when facing uncertain consumption. Convex functions are also commonly encountered, for example in option pricing or in (production) cost functions. The expectations operator has a number of simple and useful properties:

- If c is a constant, then $E[c] = c$. This property follows since the expectation is an integral against a probability density which integrates to unity.
- If c is a constant, then $E[cX] = cE[X]$. This property follows directly from passing the constant out of the integral in the definition of the expectation operator.
- The expectation of the sum is the sum of the expectations,

$$E\left[\sum_{i=1}^k g_i(X)\right] = \sum_{i=1}^k E[g_i(X)].$$

This property follows directly from the distributive property of multiplication.

- If a is a constant, then $E[a + X] = a + E[X]$. This property also follows from the distributive property of multiplication.
- $E[f(X)] = f(E[X])$ when $f(x)$ is affine (i.e. $f(x) = a + bx$ where a and b are constants). For general non-linear functions, it is usually the case that $E[f(X)] \neq f(E[X])$ when X is non-degenerate.
- $E[X^p] \neq E[X]^p$ except when $p = 1$ when X is non-degenerate.

These rules are used throughout financial economics when studying random variables and functions of random variables.

The expectation of a function of a multivariate random variable is similarly defined, only integrating across all dimensions.

Definition 1.39 (Expectation of a Multivariate Random Variable). Let (X_1, X_2, \dots, X_n) be a continuously distributed n -dimensional multivariate random variable with joint density function $f(x_1, x_2, \dots, x_n)$. The expectation of $Y = g(X_1, X_2, \dots, X_n)$ is defined as

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x_1, x_2, \dots, x_n) f(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n. \quad (1.24)$$

¹⁴A degenerate random variable has probability 1 on a single point, and so is not meaningfully random.

It is straight forward to see that rule that the expectation of the sum is the sum of the expectation carries over to multivariate random variables, and so

$$E \left[\sum_{i=1}^n g_i (X_1, \dots, X_n) \right] = \sum_{i=1}^n E [g_i (X_1, \dots, X_n)].$$

Additionally, taking $g_i (X_1, \dots, X_n) = X_i$, we have $E \left[\sum_{i=1}^n X_i \right] = \sum_{i=1}^n E [X_i]$.

1.4.2 Moments

Moments are expectations of particular functions of a random variable, typically $g(x) = x^s$ for $s = 1, 2, \dots$, and are often used to compare distributions or to estimate parameters.

Definition 1.40 (Noncentral Moment). The r^{th} noncentral moment of a continuous random variable X is defined

$$\mu'_r \equiv E[X^r] = \int_{-\infty}^{\infty} x^r f(x) dx \quad (1.25)$$

for $r = 1, 2, \dots$

The first non-central moment is the average, or mean, of the random variable.

Definition 1.41 (Mean). The first non-central moment of a random variable X is called the mean of X and is denoted μ .

Central moments are similarly defined, only centered around the mean.

Definition 1.42 (Central Moment). The r^{th} central moment of a random variables X is defined

$$\mu_r \equiv E[(X - \mu)^r] = \int_{-\infty}^{\infty} (x - \mu)^r f(x) dx \quad (1.26)$$

for $r = 2, 3, \dots$

Aside from the first moment, references to “moments” refer to central moments. Moments may not exist if a distribution is sufficiently heavy-tailed. However, if the r^{th} moment exists, then any moment of lower order must also exist.

Theorem 1.11 (Lesser Moment Existence). *If μ'_r exists for some r , then μ'_s exists for $s \leq r$. Moreover, for any r , μ'_r exists if and only if μ_r exists.*

Central moments are used to describe a distribution since they are invariant to changes in the mean. The second central moment is known as the variance.

Definition 1.43 (Variance). The second central moment of a random variable X , $E[(X - \mu)^2]$ is called the variance and is denoted σ^2 or equivalently $V[X]$.

The variance operator ($V[\cdot]$) also has a number of useful properties.

- If c is a constant, then $V[c] = 0$.
- If c is a constant, then $V[cX] = c^2V[X]$.
- If a is a constant, then $V[a + X] = V[X]$.
- The variance of the sum is the sum of the variances plus twice all of the covariances^a,

$$V\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n V[X_i] + 2 \sum_{j=1}^n \sum_{k=j+1}^n \text{Cov}[X_j, X_k]$$

^aSee Section 1.4.7 for more on covariances.

The variance is a measure of dispersion, although the square root of the variance, known as the standard deviation, is typically more useful.¹⁵

Definition 1.44 (Standard Deviation). The square root of the variance is known as the standard deviations and is denoted σ or equivalently $\text{std}(X)$.

The standard deviation is a more meaningful measure than the variance since its *units* are the same as the mean (and random variable). For example, suppose X is the return on the stock market next year, and that the mean of X is 8% and the standard deviation is 20% (the variance is .04). The mean and standard deviation are both measured as the percentage change in investment, and so can be directly compared, such as in the Sharpe ratio (Sharpe, 1994). Applying the properties of the expectation operator and variance operator, it is possible to define a studentized (or standardized) random variable.

Definition 1.45 (Studentization). Let X be a random variable with mean μ and variance σ^2 , then

$$Z = \frac{x - \mu}{\sigma} \quad (1.27)$$

is a studentized version of X (also known as standardized). Z has mean 0 and variance 1.

Standard deviation also provides a bound on the probability which can lie in the tail of a distribution, as shown in Chebyshev's inequality.

Theorem 1.12 (Chebyshev's Inequality). $\Pr[|x - \mu| \geq k\sigma] \leq 1/k^2$ for $k > 0$.

Chebyshev's inequality is useful in a number of contexts. One of the most useful is in establishing consistency in any an estimator which has a variance that tends to 0 as the sample size diverges.

The third central moment does not have a specific name, although it is called the skewness when standardized by the scaled variance.

¹⁵The standard deviation is occasionally confused for the standard error. While both are square roots of variances, the standard deviation refers to deviation in a random variable while the standard error is reserved for parameter estimators.

Definition 1.46 (Skewness). The third central moment, standardized by the second central moment raised to the power $3/2$,

$$\frac{\mu_3}{(\sigma^2)^{3/2}} = \frac{E[(X - E[X])^3]}{E[(X - E[X])^2]^{3/2}} = E[Z^3] \quad (1.28)$$

is defined as the skewness where Z is a studentized version of X .

The skewness is a general measure of asymmetry, and is 0 for symmetric distribution (assuming the third moment exists). The normalized fourth central moment is known as the kurtosis.

Definition 1.47 (Kurtosis). The fourth central moment, standardized by the squared second central moment,

$$\frac{\mu_4}{(\sigma^2)^2} = \frac{E[(X - E[X])^4]}{E[(X - E[X])^2]^2} = E[Z^4] \quad (1.29)$$

is defined as the kurtosis and is denoted κ where Z is a studentized version of X .

Kurtosis measures of the chance of observing a large (and absolute terms) value, and is often expressed as excess kurtosis.

Definition 1.48 (Excess Kurtosis). The kurtosis of a random variable minus the kurtosis of a normal random variable, $\kappa - 3$, is known as excess kurtosis.

Random variables with a positive excess kurtosis are often referred to as heavy-tailed.

1.4.3 Related Measures

While moments are useful in describing the properties of a random variable, other measures are also commonly encountered. The median is an alternative measure of central tendency.

Definition 1.49 (Median). Any number m satisfying $\Pr(X \leq m) = 0.5$ and $\Pr(X \geq m) = 0.5$ is known as the median of X .

The median measures the point where 50% of the distribution lies on either side (it may not be unique), and is just a particular quantile. The median has a few advantages over the mean, and in particular, it is less affected by outliers (e.g. the difference between mean and median income) and it always exists (the mean doesn't exist for very heavy-tailed distributions).

The interquartile range uses quartiles¹⁶ to provide an alternative measure of dispersion than standard deviation.

Definition 1.50 (Interquartile Range). The value $q_{.75} - q_{.25}$ is known as the interquartile range.

The mode complements the mean and median as a measure of central tendency. A mode is a local maximum of a density.

¹⁶Other *tiles* include terciles (3), quartiles (4), quintiles (5), deciles (10) and percentiles (100). In all cases the bin ends are $[(i - 1/m), i/m]$ where m is the number of bins and $i = 1, 2, \dots, m$.

Definition 1.51 (Mode). Let X be a random variable with density function $f(x)$. A point c where $f(x)$ attains a maximum is known as a mode.

Distributions can be unimodal or multimodal.

Definition 1.52 (Unimodal Distribution). Any random variable which has a single, unique mode is called unimodal.

Note that modes in a multimodal distribution do not necessarily have to have equal probability.

Definition 1.53 (Multimodal Distribution). Any random variable which has more than one mode is called multimodal.

Figure 1.9 shows a number of distributions. The distributions depicted in the top panels are all unimodal. The distributions in the bottom pane are mixtures of normals, meaning that with probability p random variables come from one normal, and with probability $1 - p$ they are drawn from the other. Both mixtures of normals are multimodal.

1.4.4 Multivariate Moments

Other moment definitions are only meaningful when studying 2 or more random variables (or an n -dimensional random variable). When applied to a vector or matrix, the expectations operator applies element-by-element. For example, if X is an n -dimensional random variable,

$$E[X] = E \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} E[X_1] \\ E[X_2] \\ \vdots \\ E[X_n] \end{bmatrix}. \quad (1.30)$$

Covariance is a measure which captures the tendency of two variables to move together in a linear sense.

Definition 1.54 (Covariance). The covariance between two random variables X and Y is defined

$$\text{Cov}[X, Y] = \sigma_{XY} = E[(X - E[X])(Y - E[Y])]. \quad (1.31)$$

Covariance can be alternatively defined using the joint product moment and the product of the means.

Theorem 1.13 (Alternative Covariance). *The covariance between two random variables X and Y can be equivalently defined*

$$\sigma_{XY} = E[XY] - E[X]E[Y]. \quad (1.32)$$

Inverting the covariance expression shows that no covariance is sufficient to ensure that the expectation of a product is the product of the expectations.

Theorem 1.14 (Zero Covariance and Expectation of Product). *If X and Y have $\sigma_{XY} = 0$, then $E[XY] = E[X]E[Y]$.*

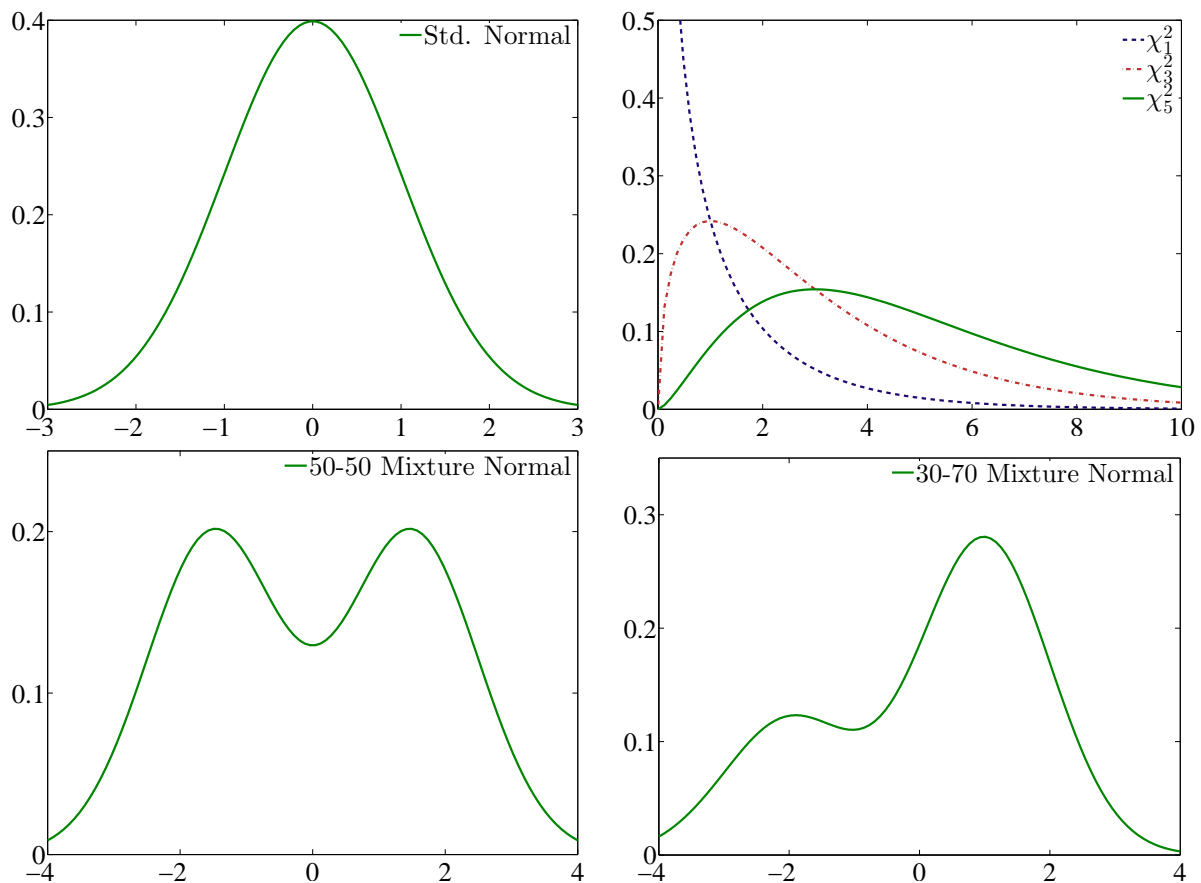


Figure 1.9: These four figures show two unimodal (upper panels) and two multimodal (lower panels) distributions. The upper left is a standard normal density. The upper right shows three χ^2 densities for $\nu = 1, 3$ and 5 . The lower panels contain mixture distributions of 2 normals – the left is a 50-50 mixture of $N(-1, 1)$ and $N(1, 1)$ and the right is a 30-70 mixture of $N(-2, 1)$ and $N(1, 1)$.

The previous result follows directly from the definition of covariance since $\sigma_{XY} = E[XY] - E[X]E[Y]$. In financial economics, this result is often applied to products of random variables so that the mean of the product can be directly determined by knowledge of the mean of each variable and the covariance between the two. For example, when studying consumption based asset pricing, it is common to encounter terms involving the expected value of consumption growth times the pricing kernel (or stochastic discount factor) – in many cases the full joint distribution of the two is intractable although the mean and covariance of the two random variables can be determined.

The Cauchy-Schwarz inequality is a version of the triangle inequality and states that the expectation of the squared product is less than the product of the squares.

Theorem 1.15 (Cauchy-Schwarz Inequality). $E[(XY)^2] \leq E[X^2] E[Y^2]$.

Example 1.32. When X is an n -dimensional random variable, it is useful to assemble the variances and covariances into a covariance matrix.

Definition 1.55 (Covariance Matrix). The covariance matrix of an n -dimensional random variable X is defined

$$\text{Cov}[X] = \Sigma = E[(X - E[X])(X - E[X])'] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \vdots & \sigma_{1n} \\ \sigma_{12} & \sigma_2^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1n} & \sigma_{2n} & \cdots & \sigma_n^2 \end{bmatrix}$$

where the i^{th} diagonal element contains the variance of X_i (σ_i^2) and the element in position (i, j) contains the covariance between X_i and X_j (σ_{ij}).

When X is composed of two sub-vectors, a block form of the covariance matrix is often convenient.

Definition 1.56 (Block Covariance Matrix). Suppose X_1 is an n_1 -dimensional random variable and X_2 is an n_2 -dimensional random variable. The block covariance matrix of $X = [X_1' X_2']'$ is

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}' & \Sigma_{22} \end{bmatrix} \quad (1.33)$$

where Σ_{11} is the n_1 by n_1 covariance of X_1 , Σ_{22} is the n_2 by n_2 covariance of X_2 and Σ_{12} is the n_1 by n_2 covariance matrix between X_1 and X_2 and element (i, j) equal to $\text{Cov}[X_{1,i}, X_{2,j}]$.

A standardized version of covariance is often used to produce a scale-free measure.

Definition 1.57 (Correlation). The correlation between two random variables X and Y is defined

$$\text{Corr}[X, Y] = \rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}. \quad (1.34)$$

Additionally, the correlation is always in the interval $[-1, 1]$, which follows from the Cauchy-Schwarz inequality.

Theorem 1.16. If X and Y are independent random variables, then $\rho_{XY} = 0$ as long as σ_X^2 and σ_Y^2 exist.

It is important to note that the converse of this statement is not true – that is, a lack of correlation does not imply that two variables are independent. In general, a correlation of 0 only implies independence when the variables are multivariate normal.

Example 1.33. Suppose X and Y have $\rho_{XY} = 0$, then X and Y are not necessarily independent. Suppose X is a discrete uniform random variable taking values in $\{-1, 0, 1\}$ and $Y = X^2$, so that $\sigma_X^2 = 2/3$, $\sigma_Y^2 = 2/9$ and $\sigma_{XY} = 0$. While X and Y are uncorrelated, they are clearly not independent, since when the random variable Y takes the value 1, X must be 0.

The corresponding correlation matrix can be assembled. Note that a correlation matrix has 1s on the diagonal and values bounded by $[-1, 1]$ on the off-diagonal positions.

Definition 1.58 (Correlation Matrix). The correlation matrix of an n -dimensional random variable X is defined

$$(\Sigma \odot \mathbf{I}_n)^{-\frac{1}{2}} \Sigma (\Sigma \odot \mathbf{I}_n)^{-\frac{1}{2}} \quad (1.35)$$

where the i, j^{th} element has the form $\sigma_{X_i X_j} / (\sigma_{X_i} \sigma_{X_j})$ when $i \neq j$ and 1 when $i = j$.

1.4.5 Conditional Expectations

Conditional expectations are similar to other forms of expectations only using conditional densities in place of joint or marginal densities. Conditional expectations essentially treat one of the variables (in a bivariate random variable) as constant.

Definition 1.59 (Bivariate Conditional Expectation). Let X be a continuous bivariate random variable comprised of X_1 and X_2 . The conditional expectation of X_1 given X_2

$$\mathbb{E} [g(X_1) | X_2 = x_2] = \int_{-\infty}^{\infty} g(x_1) f(x_1 | x_2) dx_1 \quad (1.36)$$

where $f(x_1 | x_2)$ is the conditional probability density function of X_1 given X_2 .¹⁷

In many cases, it is useful to avoid specifying a specific value for X_2 in which case $\mathbb{E} [X_1 | X_1]$ will be used. Note that $\mathbb{E} [X_1 | X_2]$ will typically be a function of the random variable X_2 .

Example 1.34. Suppose X is a bivariate normal distribution with components X_1 and X_2 , $\mu = [\mu_1 \mu_2]'$ and

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix},$$

then $\mathbb{E} [X_1 | X_2 = x_2] = \mu_1 + \frac{\sigma_{12}}{\sigma_2^2} (x_2 - \mu_2)$. This follows from the conditional density of a bivariate random variable.

The law of iterated expectations uses conditional expectations to show that the conditioning does not affect the final result of taking expectations – in other words, the order of taking expectations does not matter.

¹⁷A conditional expectation can also be defined in a natural way for functions of X_1 given $X_2 \in B$ where $\Pr(X_2 \in B) > 0$.

Theorem 1.17 (Bivariate Law of Iterated Expectations). *Let X be a continuous bivariate random variable comprised of X_1 and X_2 . Then $E[E[g(X_1)|X_2]] = E[g(X_1)]$.*

The law of iterated expectations follows from basic properties of an integral since the order of integration does not matter as long as all integrals are taken.

Example 1.35. Suppose X is a bivariate normal distribution with components X_1 and X_2 , $\mu = [\mu_1 \mu_2]'$ and

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix},$$

then $E[X_1] = \mu_1$ and

$$\begin{aligned} E[E[X_1|X_2]] &= E\left[\mu_1 + \frac{\sigma_{12}}{\sigma_2^2}(X_2 - \mu_2)\right] \\ &= \mu_1 + \frac{\sigma_{12}}{\sigma_2^2}(E[X_2] - \mu_2) \\ &= \mu_1 + \frac{\sigma_{12}}{\sigma_2^2}(\mu_2 - \mu_2) \\ &= \mu_1. \end{aligned}$$

When using conditional expectations, any random variable conditioned on behaves “as-if” non-random (in the conditional expectation), and so $E[E[X_1X_2|X_2]] = E[X_2E[X_1|X_2]]$. This is a very useful tool when combined with the law of iterated expectations when $E[X_1|X_2]$ is a known function of X_2 .

Example 1.36. Suppose X is a bivariate normal distribution with components X_1 and X_2 , $\mu = \mathbf{0}$ and

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix},$$

then

$$\begin{aligned} E[X_1X_2] &= E[E[X_1X_2|X_2]] \\ &= E[X_2E[X_1|X_2]] \\ &= E\left[X_2\left(\frac{\sigma_{12}}{\sigma_2^2}X_2\right)\right] \\ &= \frac{\sigma_{12}}{\sigma_2^2}E[X_2^2] \\ &= \frac{\sigma_{12}}{\sigma_2^2}(\sigma_2^2) \\ &= \sigma_{12}. \end{aligned}$$

One particularly useful application of conditional expectations occurs when the conditional expectation is known and constant, so that $E[X_1|X_2] = c$.

Example 1.37. Suppose X is a bivariate random variable composed of X_1 and X_2 and that $E[X_1|X_2] = c$. Then $E[X_1] = c$ since

$$\begin{aligned} E[X_1] &= E[E[X_1|X_2]] \\ &= E[c] \\ &= c. \end{aligned}$$

Conditional expectations can be taken for general n -dimensional random variables, and the law of iterated expectations holds as well.

Definition 1.60 (Conditional Expectation). Let X be a n -dimensional random variable and partition the first j ($1 \leq j < n$) elements of X into X_1 , and the remainder into X_2 so that $X = [X_1' X_2']'$. The conditional expectation of $g(X_1)$ given $X_2 = \mathbf{x}_2$

$$E[g(X_1)|X_2 = \mathbf{x}_2] = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x_1, \dots, x_j) f(x_1, \dots, x_j|\mathbf{x}_2) dx_j \dots dx_1 \quad (1.37)$$

where $f(x_1, \dots, x_j|\mathbf{x}_2)$ is the conditional probability density function of X_1 given $X_2 = \mathbf{x}_2$.

The law of iterated expectations also holds for arbitrary partitions as well.

Theorem 1.18 (Law of Iterated Expectations). Let X be a n -dimensional random variable and partition the first j ($1 \leq j < n$) elements of X into X_1 , and the remainder into X_2 so that $X = [X_1' X_2']'$. Then $E[E[g(X_1)|X_2]] = E[g(X_1)]$. The law of iterated expectations is also known as the law of total expectations.

Full multivariate conditional expectations are extremely common in time series. For example, when using daily data, there are over 30,000 observations of the Dow Jones Industrial Average available to model. Attempting to model the full joint distribution would be a formidable task. On the other hand, modeling the conditional expectation (or conditional mean) of the final observation, conditioning on those observations in the past, is far simpler.

Example 1.38. Suppose $\{X_t\}$ is a sequence of random variables where X_t comes after X_{t-j} for $j \geq 1$. The conditional conditional expectation of X_t given its past is

$$E[X_t|X_{t-1}, X_{t-2}, \dots].$$

Example 1.39. Let $\{\epsilon_t\}$ be a sequence of independent, identically distributed random variables with mean 0 and variance $\sigma^2 < \infty$. Define $X_0 = 0$ and $X_t = X_{t-1} + \epsilon_t$. X_t is a random walk, and $E[X_t|X_{t-1}] = X_{t-1}$.

This leads naturally to the definition of a martingale, which is an important concept in financial economics which related to efficient markets.

Definition 1.61 (Martingale). If $E[X_{t+j}|X_{t-1}, X_{t-2}, \dots] = X_{t-1}$ for all $j \geq 0$ and $E[|X_t|] < \infty$, both holding for all t , then $\{X_t\}$ is a martingale. Similarly, if $E[X_{t+j} - X_{t-1}|X_{t-1}, X_{t-2}, \dots] = 0$ for all $j \geq 0$ and $E[|X_t|] < \infty$, both holding for all t , then $\{X_t\}$ is a martingale.

1.4.6 Conditional Moments

All moments can be transformed made conditional by integrating against the conditional probability density function. For example, the (unconditional) mean becomes the conditional mean, and the variance becomes a conditional variance.

Definition 1.62 (Conditional Variance). The variance of a random variable X conditional on another random variable Y is

$$\begin{aligned} V[X|Y] &= E \left[(X - E[X|Y])^2 | Y \right] \\ &= E[X^2|Y] - E[X|Y]^2. \end{aligned} \quad (1.38)$$

The two definitions of conditional variance are identical to those of the (unconditional) variance where the (unconditional) expectation has been replaced by a conditional expectation. Conditioning can be used to compute higher-order moments as well.

Definition 1.63 (Conditional Moment). The r^{th} central moment of a random variables X conditional on another random variable Y is defined

$$\mu_r \equiv E \left[(X - E[X|Y])^r | Y \right] \quad (1.39)$$

for $r = 2, 3, \dots$

Combining the conditional expectation and the conditional variance leads to the law of total variance.

Theorem 1.19. *The variance of a random variable X can be decomposed into the variance of the conditional expectation plus the expectation of the conditional variance,*

$$V[X] = V[E[X|Y]] + E[V[X|Y]]. \quad (1.40)$$

The law of total variance shows that the total variance of a variable can be decomposed into the variability of the conditional mean plus the average of the conditional variance. This is a useful decomposition for time-series.

Independence can also be defined conditionally.

Definition 1.64 (Conditional Independence). Two random variables X_1 and X_2 are conditionally independent, conditional on Y , if

$$f(x_1, x_2|y) = f_1(x_1|y) f_2(x_2|y).$$

Note that random variables that are conditionally independent are not necessarily unconditionally independent.

Example 1.40. Suppose X is a trivariate normal random variable with mean $\mathbf{0}$ and covariance

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & \sigma_3^2 \end{bmatrix}$$

and define $Y_1 = x_1 + x_3$ and $Y_2 = x_2 + x_3$. Then Y_1 and Y_2 are correlated bivariate normal with mean $\mathbf{0}$ and covariance

$$\Sigma_Y = \begin{bmatrix} \sigma_1^2 + \sigma_3^2 & \sigma_3^2 \\ \sigma_3^2 & \sigma_2^2 + \sigma_3^2 \end{bmatrix},$$

but the joint distribution of Y_1 and Y_2 given X_3 is bivariate normal with mean $\mathbf{0}$ and

$$\Sigma_{Y|X_3} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

and so Y_1 and Y_2 are independent conditional on X_3 .

Other properties of unconditionally independent random variables continue to hold for conditionally independent random variables. For example, when X_1 and X_2 are independent conditional on X_3 , then the conditional covariance between X_1 and X_2 is 0 (as is the conditional correlation), and $E[E[X_1 X_2 | X_3]] = E[E[X_1 | X_3] E[X_2 | X_3]]$ – that is, the conditional expectation of the product is the product of the conditional expectations.

1.4.7 Vector and Matrix Forms

Vector and matrix forms are particularly useful in finance since portfolios are often of interest where the underlying random variables are the individual assets and the combination vector is the vector of portfolio weights.

Theorem 1.20. Let $Y = \sum_{i=1}^n c_i X_i$ where $c_i, i = 1, \dots, n$ are constants. Then $E[Y] = \sum_{i=1}^n c_i E[X_i]$. In matrix notation, $Y = \mathbf{c}'\mathbf{x}$ where \mathbf{c} is an n by 1 vector and $E[Y] = \mathbf{c}'E[X]$.

The variance of the sum is the weighted sum of the variance plus all of the covariances.

Theorem 1.21. Let $Y = \sum_{i=1}^n c_i X_i$ where c_i are constants. Then

$$V[Y] = \sum_{i=1}^n c_i^2 V[X_i] + 2 \sum_{j=1}^n \sum_{k=j+1}^n c_j c_k \text{Cov}[X_i, X_j] \quad (1.41)$$

or equivalently

$$\sigma_Y^2 = \sum_{i=1}^n c_i^2 \sigma_{X_i}^2 + 2 \sum_{j=1}^n \sum_{k=j+1}^n c_j c_k \sigma_{X_j X_k}.$$

This result can be equivalently expressed in vector-matrix notation.

Theorem 1.22. Let \mathbf{c} in an n by 1 vector and let X by an n -dimensional random variable with covariance Σ . Define $Y = \mathbf{c}'\mathbf{x}$. The variance of Y is $\sigma_Y^2 = \mathbf{c}'\text{Cov}[X]\mathbf{c} = \mathbf{c}'\Sigma\mathbf{c}$.

Note that the result holds when \mathbf{c} is replaced by a matrix \mathbf{C} .

Theorem 1.23. Let \mathbf{C} be an n by m matrix and let X be an n -dimensional random variable with mean μ_X and covariance Σ_X . Define $Y = \mathbf{C}'\mathbf{x}$. The expected value of Y is $E[Y] = \mu_Y = \mathbf{C}'E[X] = \mathbf{C}'\mu_X$ and the covariance of Y is $\Sigma_Y = \mathbf{C}'\text{Cov}[X]\mathbf{C} = \mathbf{C}'\Sigma_X\mathbf{C}$.

Definition 1.65 (Multivariate Studentization). Let X be an n -dimensional random variable with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, then

$$Z = \boldsymbol{\Sigma}^{-\frac{1}{2}} (\mathbf{x} - \boldsymbol{\mu}) \quad (1.42)$$

is a studentized version of X where $\boldsymbol{\Sigma}^{\frac{1}{2}}$ is a matrix square root such as the Cholesky factor or one based on the spectral decomposition of $\boldsymbol{\Sigma}$. Z has mean $\mathbf{0}$ and covariance equal to the identity matrix \mathbf{I}_n .

The final result for vectors relates quadratic forms of normals (inner-products) to χ^2 distributed random variables.

Theorem 1.24 (Quadratic Forms of Normals). Let X be an n -dimensional normal random variable with mean $\mathbf{0}$ and identity covariance \mathbf{I}_n . Then $\mathbf{x}'\mathbf{x} = \sum_{i=1}^n x_i^2 \sim \chi_n^2$.

Combing this result with studentization, when X is a general n -dimensional normal random variable with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$,

$$(\mathbf{x} - \boldsymbol{\mu})' \left(\boldsymbol{\Sigma}^{-\frac{1}{2}} \right)' \boldsymbol{\Sigma}^{-\frac{1}{2}} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \sim \chi_n^2.$$

1.4.8 Monte Carlo and Numerical Integration

Expectations of functions of continuous random variables are integrals against the underlying pdf. In some cases, these integrals are analytically tractable, although in many situations integrals cannot be analytically computed and so numerical techniques are needed to compute expected values and moments.

Monte Carlo is one method to approximate an integral. Monte Carlo utilizes simulated draws from the underlying distribution and averaging to approximate integrals.

Definition 1.66 (Monte Carlo Integration). Suppose $X \sim F(\theta)$ and that it is possible to simulate a series $\{x_i\}$ from $F(\theta)$. The Monte Carlo expectation of a function $g(x)$ is defined

$$E[\widehat{g}(X)] = m^{-1} \sum_{i=1}^m g(x_i),$$

Moreover, as long as $E[|g(x)|] < \infty$, $\lim_{m \rightarrow \infty} m^{-1} \sum_{i=1}^m g(x_i) = E[g(x)]$.

The intuition behind this result follows from the properties of $\{x_i\}$. Since these are i.i.d. draws from $F(\theta)$, they will, on average, tend to appear in any interval $B \in R(X)$ in proportion to the probability $\Pr(X \in B)$. In essence, the simulated values coarsely approximating the discrete approximation shown in 1.8.

While Monte Carlo integration is a general technique, there are some important limitations. First, if the function $g(x)$ takes large values in regions where $\Pr(X \in B)$ is small, it may require a very large number of draws to accurately approximate $E[g(x)]$ since, by construction, there are unlikely to many points in B . In practice the behavior of $h(x) = g(x)f(x)$ plays an important role in determining the appropriate sample size.¹⁸ Second, while Monte Carlo integration

¹⁸Monte Carlo integrals can also be seen as estimators, and in many cases standard inference can be used to determine the accuracy of the integral. See Chapter 1 for more details on inference and constructing confidence intervals.

is technically valid for random variables with any number of dimensions, in practice it is usually only reliable when the dimension is small (typically 3 or fewer), especially when the range is unbounded ($R(X) \in \mathbb{R}^n$). When the dimension of X is large, many simulated draws are needed to visit the corners of the (joint) pdf, and if 1,000 draws are sufficient for a unidimensional problem, 1000^n may be needed to achieve the same accuracy when X has n dimensions.

Alternatively the function to be integrated can be approximated using a polygon with an easy-to-compute area, such as the rectangles approximating the normal pdf shown in figure 1.8. The quality of the approximation will depend on the resolution of the grid used. Suppose u and l are the upper and lower bounds of the integral, respectively, and that the region can be split into m intervals $l = b_0 < b_1 < \dots < b_{m-1} < b_m = u$. Then the integral of a function $h(\cdot)$ is

$$\int_l^u h(x) dx = \sum_{i=1}^m \int_{b_{i-1}}^{b_i} h(x) dx.$$

In practice, l and u may be infinite, in which case some cut-off point is required. In general, the cut-off should be chosen so that the vast majority of the probability lies between l and u ($\int_l^u f(x) dx \approx 1$).

This decomposition is combined with an area for approximating the area under h between b_{i-1} and b_i . The simplest is the rectangle method, which uses a rectangle with a height equal to the value of the function at the mid-point.

Definition 1.67 (Rectangle Method). The rectangle rule approximates the area under the curve with a rectangle and is given by

$$\int_l^u h(x) dx \approx h\left(\frac{u+l}{2}\right)(u-l).$$

The rectangle rule would be exact if the function was piece-wise flat. The trapezoid rule improves the approximation by replacing the function at the midpoint with the average value of the function and would be exact for any piece-wise linear function (including piece-wise flat functions).

Definition 1.68 (Trapezoid Method). The trapezoid rule approximates the area under the curve with a trapezoid and is given by

$$\int_l^u h(x) dx \approx \frac{h(u) + h(l)}{2}(u-l).$$

The final method is known as Simpson's rule which is based on using a quadratic approximation to the underlying function. It is exact when the underlying function is piece-wise linear or quadratic.

Definition 1.69 (Simpson's Rule). Simpson's Rule uses an approximation that would be exact if the underlying function were quadratic, and is given by

$$\int_l^u h(x) dx \approx \frac{u-l}{6} \left(h(u) + 4h\left(\frac{u+l}{2}\right) + h(l) \right).$$

Example 1.41. Consider the problem of computing the expected payoff of an option. The payoff of a call option is given by

$$c = \max(s_1 - k, 0)$$

where k is the strike price, s_1 is the stock price at expiration and s_0 is the current stock price. Suppose returns are normally distributed with mean $\mu = .08$ and standard deviation $\sigma = .20$. In this problem, $g(r) = (s_0 \exp(r) - k) I_{[s_0 \exp(r) > k]}$ where $I_{[\cdot]}$ is a binary indicator function which takes the value 1 when the argument is true, and

$$f(r) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(r - \mu)^2}{2\sigma^2}\right).$$

Combined, the function to be integrated is

$$\begin{aligned} \int_{-\infty}^{\infty} h(r) dr &= \int_{-\infty}^{\infty} g(r) f(r) dr \\ &= \int_{-\infty}^{\infty} (s_0 \exp(r) - k) I_{[s_0 \exp(r) > k]} \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(r - \mu)^2}{2\sigma^2}\right) dr \end{aligned}$$

$s_0 = k = 50$ was used in all results.

All four methods were applied to the problem. The number of bins and the range of integration was varied for the analytical approximations. The number of bins ranged across $\{10, 20, 50, 1000\}$ and the integration range spanned $\{\pm 3\sigma, \pm 4\sigma, \pm 6\sigma, \pm 10\sigma\}$ and the bins were uniformly spaced along the integration range. Monte Carlo integration was also applied with $m \in \{100, 1000\}$.

All things equal, increasing the number of bins increases the accuracy of the approximation. In this example, 50 appears to be sufficient. However, having a range which is too small produces values which differ from the correct value of 7.33. The sophistication of the method also improves the accuracy, especially when the number of nodes is small. The Monte Carlo results are also close, on average. However, the standard deviation is large, about 5%, even when 1000 draws are used, so that large errors would be commonly encountered and so many more points are needed to ensure that the integral is always accurate.

Shorter Problems

Problem 1.1. Suppose

$$\begin{bmatrix} X \\ U \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_X \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_X^2 & 0 \\ 0 & \sigma_U^2 \end{bmatrix}\right)$$

and $Y = 2X + U$. What is $E[Y]$ and $V[Y]$?

Problem 1.2. Show $\text{Cov}[aX + bY, cX + dY] = acV[X] + bdV[Y] + (ad + bc)\text{Cov}[X, Y]$.

Problem 1.3. Show that the two forms of the covariance,

$$E[XY] - E[X]E[Y] \text{ and } E[(X - E[X])(Y - E[Y])]$$

are equivalent when X and Y are continuous random variables.

Problem 1.4. Suppose $\{X_i\}$ is a sequence of random variables where $V[X_i] = \sigma^2$ for all i , $\text{Cov}[X_i, X_{i-1}] = \theta$ and $\text{Cov}[X_i, X_{i-j}] = 0$ for $j > 1$. What is $V[\bar{X}]$ where $\bar{X} = n^{-1} \sum_{i=1}^n X_i$?

Problem 1.5. Suppose $Y = \beta X + \epsilon$ where $X \sim N(\mu_X, \sigma_X^2)$, $\epsilon \sim N(0, \sigma^2)$ and X and ϵ are independent. What is $\text{Corr}[X, Y]$?

Problem 1.6. Prove that $E[a + bX] = a + bE[X]$ when X is a continuous random variable.

Problem 1.7. Prove that $V[a + bX] = b^2V[X]$ when X is a continuous random variable.

Problem 1.8. Prove that $\text{Cov}[a + bX, c + dY] = bd\text{Cov}[X, Y]$ when X and Y are a continuous random variables.

Problem 1.9. Prove that $V[a + bX + cY] = b^2V[X] + c^2V[Y] + 2bc\text{Cov}[X, Y]$ when X and Y are a continuous random variables.

Problem 1.10. Suppose $\{X_i\}$ is an i.i.d. sequence of random variables. Show that

$$V[\bar{X}] = V\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = n^{-1}\sigma^2$$

where σ^2 is $V[X_1]$.

Problem 1.11. Prove that $\text{Corr}[a + bX, c + dY] = \text{Corr}[X, Y]$.

Problem 1.12. Suppose $\{X_i\}$ is a sequence of random variables where, for all i , $V[X_i] = \sigma^2$, $\text{Cov}[X_i, X_{i-1}] = \theta$ and $\text{Cov}[X_i, X_{i-j}] = 0$ for $j > 1$. What is $V[\bar{X}]$?

Problem 1.13. Prove that $E[a + bX|Y] = a + bE[X|Y]$ when X and Y are continuous random variables.

Problem 1.14. Suppose that $E[X|Y] = Y^2$ where Y is normally distributed with mean μ and variance σ^2 . What is $E[a + bX]$?

Problem 1.15. Suppose $E[X|Y = y] = a + by$ and $V[X|Y = y] = c + dy^2$ where Y is normally distributed with mean μ and variance σ^2 . What is $V[X]$?

Problem 1.16. Show that the law of total variance holds for a $V[X_1]$ when X is a bivariate normal with mean $\mu = [\mu_1 \mu_2]'$ and covariance

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}.$$

Longer Exercises

Exercise 1.1. Sixty percent (60%) of all traders hired by a large financial firm are rated as performing satisfactorily or better in their first-year review. Of these, 90% earned a first in financial econometrics. Of the traders who were rated as unsatisfactory, only 20% earned a first in financial econometrics.

1. What is the probability that a trader is rated as satisfactory or better given they received a first in financial econometrics?
2. What is the probability that a trader is rated as unsatisfactory given they received a first in financial econometrics?
3. Is financial econometrics a useful indicator of trader performance? Why or why not?

Exercise 1.2. Large financial firms use automated screening to detect rogue trades – those that exceed risk limits. One of your colleagues has introduced a new statistical test using the trading data that, given that a trader has exceeded her risk limit, detects this with probability 98%. It also only indicates false positives – that is non-rogue trades that are flagged as rogue – 1% of the time.

1. Assuming 99% of trades are legitimate, what is the probability that a detected trade is rogue? Explain the intuition behind this result.
2. Is this a useful test? Why or why not?
3. How low would the false positive rate have to be to have a 98% chance that a detected trade was actually rogue?

Exercise 1.3. Your corporate finance professor uses a few jokes to add levity to his lectures. Each week he tells 3 different jokes. However, he is also very busy, and so forgets week to week which jokes were used.

1. Assuming he has 12 jokes, what is the probability of 1 repeat across 2 consecutive weeks?
2. What is the probability of hearing 2 of the same jokes in consecutive weeks?
3. What is the probability that all 3 jokes are the same?
4. Assuming the term is 8 weeks long, and they your professor has 96 jokes, what is the probability that there is no repetition during the term? Note that he remembers the jokes he gives in a particular lecture, only forgets across lectures.
5. How many jokes would your professor need to know to have a 99% chance of not repeating any in the term?

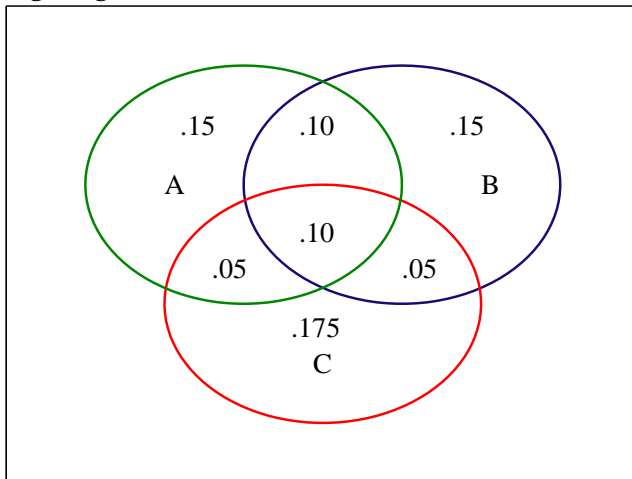
Exercise 1.4. A hedge fund company manages three distinct funds. In any given month, the probability that the return is positive is shown in the following table:

$$\begin{aligned} \Pr(r_{1,t} > 0) &= .55 & \Pr(r_{1,t} > 0 \cup r_{2,t} > 0) &= .82 \\ \Pr(r_{2,t} > 0) &= .60 & \Pr(r_{1,t} > 0 \cup r_{3,t} > 0) &= .7525 \\ \Pr(r_{3,t} > 0) &= .45 & \Pr(r_{2,t} > 0 \cup r_{3,t} > 0) &= .78 \\ \Pr(r_{2,t} > 0 \cap r_{3,t} > 0 | r_{1,t} > 0) &= .20 \end{aligned}$$

1. Are the events of “positive returns” pairwise independent?
2. Are the events of “positive returns” independent?

3. What is the probability that funds 1 and 2 have positive returns, given that fund 3 has a positive return?
4. What is the probability that at least one fund will have a positive return in any given month?

Exercise 1.5. Suppose the probabilities of three events, A , B and C are as depicted in the following diagram:



1. Are the three events pairwise independent?
2. Are the three events independent?
3. What is $\Pr(A \cap B)$?
4. What is $\Pr(A \cap B | C)$?
5. What is $\Pr(C | A \cap B)$?
6. What is $\Pr(C | A \cup B)$?

Exercise 1.6. At a small high-frequency hedge fund, two competing algorithms produce trades. Algorithm α produces 80 trades per second and 5% lose money. Algorithm β produces 20 trades per second but only 1% lose money. Given the last trade lost money, what is the probability it was produced by algorithm β ?

Exercise 1.7. Suppose $f(x, y) = 2 - x - y$ where $x \in [0, 1]$ and $y \in [0, 1]$.

1. What is $\Pr(X > .75 \cap Y > .75)$?
2. What is $\Pr(X + Y > 1.5)$?
3. Show formally whether X and Y are independent.
4. What is $\Pr(Y < .5 | X = x)$?

Exercise 1.8. Suppose $f(x, y) = xy$ for $x \in [0, 1]$ and $y \in [0, 2]$.

1. What is the joint cdf?
2. What is $\Pr(X < 0.5 \cap Y < 1)$?
3. What is the marginal cdf of X ? What is $\Pr(X < 0.5)$?
4. What is the marginal density of X ?
5. Are X and Y independent?

Exercise 1.9. Suppose $F(x) = 1 - p^{x+1}$ for $x \in [0, 1, 2, \dots]$ and $p \in (0, 1)$.

1. Find the pmf.
2. Verify that the pmf is valid.
3. What is $\Pr(X \leq 8)$ if $p = .75$?
4. What is $\Pr(X \leq 1)$ given $X \leq 8$?

Exercise 1.10. A firm producing widgets has a production function $q(L) = L^{0.5}$ where L is the amount of labor. Sales prices fluctuate randomly and can be \$10 (20%), \$20 (50%) or \$30 (30%). Labor prices also vary and can be \$1 (40%), 2 (30%) or 3 (30%). The firm always maximizes profits after seeing both sales prices and labor prices.

1. Define the distribution of profits possible?
2. What is the probability that the firm makes at least \$100?
3. Given the firm makes a profit of \$100, what is the probability that the profit is over \$200?

Exercise 1.11. A fund manager tells you that her fund has non-linear returns as a function of the market and that his return is $r_{i,t} = 0.02 + 2r_{m,t} - 0.5r_{m,t}^2$ where $r_{i,t}$ is the return on the fund and $r_{m,t}$ is the return on the market.

1. She tells you her expectation of the market return this year is 10%, and that her fund will have an expected return of 22%. Can this be?
2. At what variance is would the expected return on the fund be negative?

Exercise 1.12. For the following densities, find the mean (if it exists), variance (if it exists), median and mode, and indicate whether the density is symmetric.

1. $f(x) = 3x^2$ for $x \in [0, 1]$
2. $f(x) = 2x^{-3}$ for $x \in [1, \infty)$
3. $f(x) = [\pi(1 + x^2)]^{-1}$ for $x \in (-\infty, \infty)$
4. $f(x) = \binom{4}{x} .2^x .8^{4-x}$ for $x \in \{0, 1, 2, 3, 4\}$

Exercise 1.13. The daily price of a stock has an average value of £2. Then then $\Pr(X > 10) < .2$ where X denotes the price of the stock. True or false?

Exercise 1.14. An investor can invest in stocks or bonds which have expected returns and co-variances as

$$\mu = \begin{bmatrix} .10 \\ .03 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} .04 & -.003 \\ -.003 & .0009 \end{bmatrix}$$

where stocks are the first component.

1. Suppose the investor has £1,000 to invest and splits the investment evenly. What is the expected return, standard deviation, variance and Sharpe Ratio (μ/σ) for the investment?
2. Now suppose the investor seeks to maximize her expected utility where her utility is defined in terms of her portfolio return, $U(r) = E[r] - .01V[r]$. How much should she invest in each asset?

Exercise 1.15. Suppose $f(x) = (1 - p)^x p$ for $x \in (0, 1, \dots)$ and $p \in (0, 1]$. Show that a random variable from the distribution is “memoryless” in the sense that $\Pr(X \geq s + r | X \geq r) = \Pr(X \geq s)$. In other words, the probability of surviving s or more periods is the same whether starting at 0 or after having survived r periods.

Exercise 1.16. Your Economics professor offers to play a game with you. You pay £1,000 to play and your Economics professor will flip a fair coin and pay you 2^x where x is the number of tries required for the coin to show heads.

1. What is the pmf of X ?
2. What is the expected payout from this game?

Exercise 1.17. Consider the roll of a fair pair of dice where a roll of a 7 or 11 pays $2x$ and anything else pays $-x$ where x is the amount bet. Is this game fair?

Exercise 1.18. Suppose the joint density function of X and Y is given by $f(x, y) = \frac{1}{2}x \exp(-xy)$ where $x \in [3, 5]$ and $y \in (0, \infty)$.

1. Give the form of $E[Y | X = x]$.
2. Graph the conditional expectation curve.

Exercise 1.19. Suppose a fund manager has \$10,000 of yours under management and tells you that the expected value of your portfolio in two years time is \$30,000 and that with probability 75% your investment will be worth at least \$40,000 in two years time.

1. Do you believe her?
2. Next, suppose she tells you that the standard deviation of your portfolio value is 2,000. Assuming this is true (as is the expected value), what is the most you can say about the probability your portfolio value falls between \$20,000 and \$40,000 in two years time?

Exercise 1.20. Suppose the joint probability density function of two random variables is given by $f(x, y) = \frac{2}{5}(3x + 2y)$ where $x \in [0, 1]$ and $y \in [0, 1]$.

1. What is the marginal probability density function of X ?
2. What is $E[X|Y = y]$? Are X and Y independent? (Hint: What must the form of $E[X|Y]$ be when they are independent?)

Exercise 1.21. Let Y be distributed χ^2_{15} .

1. What is $\Pr(y > 27.488)$?
2. What is $\Pr(6.262 \leq y \leq 27.488)$?
3. Find C where $\Pr(y \geq c) = \alpha$ for $\alpha \in \{0.01, 0.05, 0.01\}$.
Next, Suppose Z is distributed χ^2_5 and is independent of Y .
4. Find C where $\Pr(y + z \geq c) = \alpha$ for $\alpha \in \{0.01, 0.05, 0.01\}$.

Exercise 1.22. Suppose X is a bivariate random variable with parameters

$$\mu = \begin{bmatrix} 5 \\ 8 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix}.$$

1. What is $E[X_1|X_2]$?
2. What is $V[X_1|X_2]$?
3. Show (numerically) that the law of total variance holds for X_2 .

Exercise 1.23. Suppose $y \sim N(5, 36)$ and $x \sim N(4, 25)$ where X and Y are independent.

1. What is $\Pr(y > 10)$?
2. What is $\Pr(-10 < y < 10)$?
3. What is $\Pr(x - y > 0)$?
4. Find C where $\Pr(x - y > C) = \alpha$ for $\alpha \in \{0.10, 0.05, 0.01\}$?

Rectangle Method				
Bins	$\pm 3\sigma$	$\pm 4\sigma$	$\pm 6\sigma$	$\pm 10\sigma$
10	7.19	7.43	7.58	8.50
20	7.13	7.35	7.39	7.50
50	7.12	7.33	7.34	7.36
1000	7.11	7.32	7.33	7.33

Trapezoid Method				
Bins	$\pm 3\sigma$	$\pm 4\sigma$	$\pm 6\sigma$	$\pm 10\sigma$
10	6.96	7.11	6.86	5.53
20	7.08	7.27	7.22	7.01
50	7.11	7.31	7.31	7.28
1000	7.11	7.32	7.33	7.33

Simpson's Rule				
Bins	$\pm 3\sigma$	$\pm 4\sigma$	$\pm 6\sigma$	$\pm 10\sigma$
10	7.11	7.32	7.34	7.51
20	7.11	7.32	7.33	7.34
50	7.11	7.32	7.33	7.33
1000	7.11	7.32	7.33	7.33

Monte Carlo		
Draws (m)	100	1000
Mean	7.34	7.33
Std. Dev.	0.88	0.28

Table 1.1: Computed values for the expected payout for an option, where the correct value is 7.33. The top three panels use approximations to the function which have simple to compute areas. The bottom panel shows the average and standard deviation from a Monte Carlo integration where the number of points varies and 10,000 simulations were used.