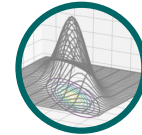


Probability and Distributions



Probability, loosely speaking, concerns the study of uncertainty. Probability can be thought of as the fraction of times an event occurs, or as a degree of belief about an event. We then would like to use this probability to measure the chance of something occurring in an experiment. As mentioned in Chapter 1, we often quantify: uncertainty in the data, uncertainty in the machine learning model, and uncertainty in the predictions produced by the model. Quantifying uncertainty requires the idea of a *random variable*, which is a function that maps outcomes of random experiments to a set of properties that we are interested in. Associated with the random variable is a function which measures the probability that a particular outcome (or set of outcomes) will occur, this is called the *probability distribution*.

random variable

probability
distribution

Probability distributions are used as a building block for other concepts, such as probabilistic modeling (Section 8.3), graphical models (Section 8.4) and model selection (Section 8.5). In the next section, we present the three concepts that define a probability space (the state space, the events and the probability of an event) and how they are related to a fourth concept called the random variable. The presentation is deliberately slightly hand wavy since a rigorous presentation would occlude the main idea. An outline of the concepts presented in this chapter are shown in Figure 6.1.

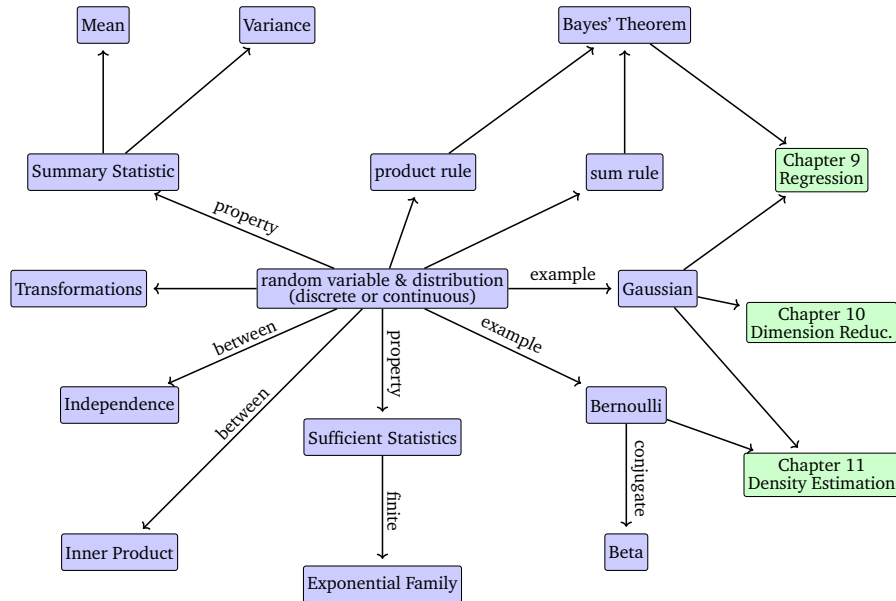
6.1 Construction of a Probability Space

The theory of probability aims at defining a mathematical structure to describe random outcomes of experiments. For example, when tossing a single coin, one cannot determine the outcome, but by doing a large number of coin tosses, one can observe a regularity in the average outcome. Using this mathematical structure of probability, the goal is to perform automated reasoning, and in this sense probability generalizes logical reasoning (Jaynes, 2003).

6.1.1 Philosophical Issues

When constructing automated reasoning systems, classical Boolean logic does not allow us to express certain forms of plausible reasoning. Consider

Figure 6.1 A mind map of the concepts related to random variables and probability distributions, as described in this chapter.



“For plausible reasoning it is necessary to extend the discrete true and false values of truth to continuous plausibilities.” (Jaynes, 2003)

the following scenario: We observe that A is false. We find B becomes less plausible although no conclusion can be drawn from classical logic. We observe that B is true. It seems A becomes more plausible. We use this form of reasoning daily. We are waiting for a friend, and consider three possibilities; H1: she is on time, H2: she has been delayed by traffic and H3: she has been abducted by aliens. When we observe our friend is late, we must logically rule out H1. We also tend to consider H2 to be more likely, though we are not logically required to do so. Finally, we may consider H3 to be possible, but we continue to consider it quite unlikely. How do we conclude H2 is the most plausible answer? Seen in this way, probability theory can be considered a generalization of Boolean logic. In the context of machine learning, it is often applied in this way to formalize the design of automated reasoning systems. Further arguments about how probability theory is the foundation of reasoning systems can be found in (Pearl, 1988).

The philosophical basis of probability and how it should be somehow related to what we think should be true (in the logical sense) was studied by Cox (Jaynes, 2003). Another way to think about it is that if we are precise about our common sense we end up constructing probabilities. E.T. Jaynes (1922–1998) identified three mathematical criteria, which must apply to all plausibilities:

- 1 The degrees of plausibility are represented by real numbers.
- 2 These numbers must be based on the rules of common sense.
 - a) Consistency or non-contradiction: when the same result can be reached

through different means, the same plausibility value must be found in all cases.

b) Honesty: All available data must be taken into account.

c) Reproducibility: If our state of knowledge about two problems are the same, then we must assign the same degree of plausibility to both of them.

The Cox-Jaynes's theorem proves these plausibilities to be sufficient to define the universal mathematical rules that apply to plausibility p , up to transformation by an arbitrary monotonic function. Crucially, these rules are the rules of probability.

Remark. In machine learning and statistics, there are two major interpretations of probability: the Bayesian and frequentist interpretations (Bishop, 2006). The Bayesian interpretation uses probability to specify the degree of uncertainty that the user has about an event, and is sometimes referred to as subjective probability or degree of belief. The frequentist interpretation considers probability to be the relative frequencies of events, in the limit when one has infinite data. \diamond

Some machine learning texts on probabilistic models use lazy notation and jargon, which is confusing. Multiple distinct concepts are all referred to as “probability distribution”, and the reader has to often disentangle the meaning from the context. One trick to help make sense of probability distributions is to check whether we are trying to model something categorical (a discrete random variable) or something continuous (a continuous random variable). The kinds of questions we tackle in machine learning are closely related to whether we are considering categorical or continuous models.

6.1.2 Probability and Random Variables

Modern probability is based on a set of axioms proposed by Kolmogorov (Jacod and Protter, 2004, Chapter 1 and 2) that introduce the three concepts of state space, event space and probability measure. The probability space models a real world process (referred to as an experiment) with random outcomes.

The state space Ω

The *state space* is the set of all possible outcomes of the experiment, usually denoted by Ω . For example, two successive coin tosses have a state space of $\{hh, tt, ht, th\}$, where “h” denotes “heads” and “t” denotes “tails”. state space

The event space \mathcal{A}

The *event space* is the space of potential results of the experiment. A subset A of the state space Ω is in the event space \mathcal{A} if at the end of the experiment we can observe whether a particular state $\omega \in \Omega$ is in A . event space

3110 The event space \mathcal{A} is obtained by considering the collection of subsets
 3111 of Ω , and for discrete probability distributions (Section 6.2.1) \mathcal{A} is the
 3112 powerset of Ω .

3113 The probability P

probability

3114 With each event $A \in \mathcal{A}$, we associate a number $P(A)$ that measures the
 3115 probability or degree of belief that the event will occur. $P(A)$ is called
 3116 the *probability* of A .

random variable

3117 The probability of a single event must lie in the interval $[0, 1]$, and
 3118 the total probability over all states in the state space Ω must be 1, i.e.,
 3119 $P(\Omega) = 1$. Given a probability space (Ω, \mathcal{A}, P) we are interested to use it
 3120 to model some real world phenomenon. Therefore we introduce a function
 3121 $x : \Omega \rightarrow \mathcal{T}$ that takes an element of Ω (an event) and returns a
 3122 particular quantity of interest, where \mathcal{T} . For example in the case of toss-
 3123 ing two coins and counting the number of heads, a random variable x
 3124 maps to the three possible events: $x(\text{hh}) = 2$, $x(\text{ht}) = 1$, $x(\text{th}) = 1$ and
 3125 $x(\text{tt}) = 0$. In this particular case $\mathcal{T} = \{0, 1, 2\}$. This association or map-
 3126 ping from Ω to \mathcal{T} is called a *random variable*. The name “random variable”
 3127 is a great source of misunderstanding as it is neither random nor is it a
 3128 variable. It is a function. For a finite state space Ω and finite \mathcal{T} , the func-
 3129 tion corresponding to a random variable is essentially a look up table. For
 3130 any subset $S \subseteq \mathcal{T}$ we associate $P_x(S) \in [0, 1]$ (the probability) to a partic-
 3131 ular event occurring corresponding to the random variable x . Example 6.1
 3132 provides a concrete example illustrating the above terminology.

3133 *Remark.* The state space Ω above unfortunately is referred to by differ-
 3134 ent names in different books. Another common name for Ω is sample
 3135 space (Grinstead and Snell, 1997; Jaynes, 2003), and state space is some-
 3136 times reserved for referring to states in a dynamical system (Hasselblatt
 3137 and Katok, 2003). Other names sometimes used to describe Ω are: sample
 3138 description space, possibility space and (very confusingly) event space.

3139



Example 6.1

This toy example is essentially a biased coin flip example.

We assume that the reader is already familiar with computing probabilities of intersections and unions of sets of events. A more gentle introduction to probability with many examples can be found in Chapter 2 of Walpole et al. (2011).

Consider a statistical experiment where we model a funfair game consisting of drawing two coins from a bag (with replacement). There are coins from USA (denoted as \$) and UK (denoted as £) in the bag, and since we draw two coins from the bag, there are four outcomes in total. The state space or sample space Ω of this experiment is then $(\$, \$)$, $(\$, £)$, $(£, \$)$, $(£, £)$. Let us assume that the composition of the bag of coins is such that a draw returns at random a \$ with probability 0.3.

The event we are interested in is the total number of times the repeated draw returns \$. Let us define a random variable x that maps the state space Ω to \mathcal{T} , that denotes the number of times we draw \$ out of the bag. We can see from the above state space we can get zero \$, one \$ or two \$s, and therefore $\mathcal{T} = \{0, 1, 2\}$. The random variable x (a function or look up table) can be represented as a table like below

$$x((\$,\$)) = 2 \quad (6.1)$$

$$x((\$,\mathcal{L})) = 1 \quad (6.2)$$

$$x((\mathcal{L},\$)) = 1 \quad (6.3)$$

$$x((\mathcal{L},\mathcal{L})) = 0. \quad (6.4)$$

Since we return the first coin we draw before drawing the second, this implies that the two draws are independent of each other, which we will discuss in Section 6.4.5. Note that there are two experimental outcomes which map to the same event, where only one of the draws return \$. Therefore the probability mass function (Section 6.2.1) of x is given by the calculations below

$$\begin{aligned} P(x = 2) &= P((\$,\$)) \\ &= P(\$) \times P(\$) \\ &= 0.3 \times 0.3 = 0.09 \end{aligned} \quad (6.5)$$

$$\begin{aligned} P(x = 1) &= P((\$,\mathcal{L}) \cup (\mathcal{L},\$)) \\ &= P((\$,\mathcal{L})) + P((\mathcal{L},\$)) \\ &= 0.3 \times (1 - 0.3) + (1 - 0.3) \times 0.3 = 0.42 \end{aligned} \quad (6.6)$$

$$\begin{aligned} P(x = 0) &= P((\mathcal{L},\mathcal{L})) \\ &= P(\mathcal{L}) \times P(\mathcal{L}) \\ &= (1 - 0.3) \times (1 - 0.3) = 0.49. \end{aligned} \quad (6.7)$$

In the above calculation, notice that we have equated two different concepts, the probability of the output of x and the probability of the states in Ω . For example in (6.7) we say $P(x = 0) = P((\mathcal{L},\mathcal{L}))$.

Consider the random variable $x : \Omega \rightarrow \mathcal{T}$ and a subset $S \subseteq \mathcal{T}$. Let $x^{-1}(S)$ be the pre-image of S by x , that is the set of elements of Ω that map to S under x ; $\{\omega \in \Omega : x(\omega) \in S\}$. Another way to understand the transformation of probability from events in Ω via the random variable x is to associate it with the probability of the pre-image of S (Jacod and Protter, 2004). That is for $S \subseteq \mathcal{T}$,

$$P_x(S) = P(x \in S) = P(x^{-1}(S)) = P(\{\omega \in \Omega : x(\omega) \in S\}). \quad (6.8)$$

3140 The left hand side of (6.8) is the probability of the set of possible outcomes
3141 (e.g. number of heads = 1) that we are interested in. Via the random

variable x that maps states to outcomes, we see in the right hand side of (6.8) that this is the probability of the set of states (in Ω) that have the property (e.g. ht, th). We say that a random variable x is distributed according to a particular probability distribution P_x , which defines the probability mapping between the event and the probability of the outcome of the random variable. The two concepts are intertwined, but for ease of presentation we will discuss some properties with respect to random variables and others with respect to their distributions.

Remark. The range of the random variable \mathcal{T} is used to indicate the kind of probability space, that is a \mathcal{T} random variable. When \mathcal{T} is finite or countably infinite, this is called a discrete random variable (Section 6.2.1). For continuous random variables (Section 6.2.2) we only consider $\mathcal{T} = \mathbb{R}$ or $\mathcal{T} = \mathbb{R}^d$. \diamond

6.1.3 Statistics

Probability theory and statistics are often presented together, but they concern different aspects of uncertainty. One way of contrasting them is by the kinds of problems that are considered. Using probability we can consider a model of some process where the underlying uncertainty is captured by random variables, and we use the rules of probability to derive what happens. Using statistics we observe that something has happened, and try to figure out the underlying process that explains the observations. In this sense machine learning is close to statistics in its goals, that is to construct a model that adequately represents the process that generated the data. When the machine learning model is a probabilistic model, we can use the rules of probability to calculate the “best fitting” model for some data.

Another aspect of machine learning systems is that we are interested in generalization error (see Chapter 8). This means that we are actually interested in the performance of our system on instances that we will observe in future, which are not identical to the instances that we have seen so far. This analysis of future performance relies on probability and statistics, most of which is beyond what will be presented in this chapter. The interested reader is encouraged to look at the books by Shalev-Shwartz and Ben-David (2014); Boucheron et al. (2013). We will see more about statistics in Chapter 8.

6.2 Discrete and Continuous Probabilities

Let us focus our attention on ways to describe the probability of an event as introduced in Section 6.1. Depending on whether the state space is discrete or continuous the natural way to refer to distributions is different. When the state space Ω is discrete, we can specify the probability that a random variable x takes a particular value $x \in \Omega$, denoted as $P(x = x)$.

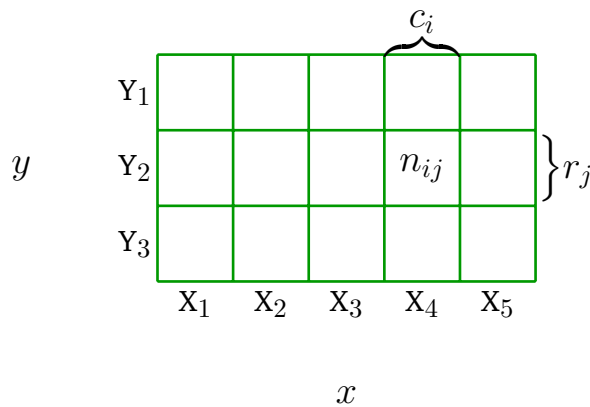


Figure 6.2
Visualization of a discrete bivariate probability mass function, with random variables x and y . This diagram is from Bishop (2006).

The expression $P(x = x)$ for a discrete random variable x is known as the *probability mass function*. When the state space Ω is continuous, e.g., the real line \mathbb{R} , it is more natural to specify the probability that a random variable x is in an interval. By convention we specify the probability that a random variable x is less than a particular value x , denoted $P(x \leq x)$. The expression $P(x \leq x)$ for a continuous random variable x is known as the *cumulative distribution function*. We will discuss continuous random variables in Section 6.2.2. We will revisit the nomenclature and contrast discrete and continuous random variables in Section 6.2.3.

Remark. We will use the phrase *univariate* distribution to refer to distributions of only one random variable (denoted by non-bold x). We will refer to distributions of more than one random variable as *multivariate* distributions, and will usually consider a vector of random variables (denoted by bold \mathbf{x}). \diamond

Many probability textbooks tend to use capital letters X for random variables and small letters x for their values. probability mass function cumulative distribution function univariate multivariate

6.2.1 Discrete Probabilities

When the state space is discrete, we can imagine the probability distribution of multiple random variables as filling out a (multidimensional) array of numbers (for example in Figure 6.2). The state space of the joint probability is the Cartesian product of the state spaces of each of the random variables. We define the *joint probability* as the entry of both values jointly

$$P(x = x_i, y = y_j) = \frac{n_{ij}}{N}, \quad (6.9)$$

where n_{ij} is the number of events with x_i and y_j and N the total number of events. The joint probability is probability of the intersection of both events, that is $P(x = x_i, y = y_j) = P(x = x_i \cap y = y_j)$. Figure 6.2 illustrates the *probability mass function* (pmf) of a discrete probability distribution. For two random variables x and y , the probability that $x = x$ and $y = y$ is (lazily) written as $p(x, y)$ and is called the joint probability.

joint probability

probability mass function

marginal probability
conditional
probability

The *marginal probability* is obtained by summing over a row or column. The *conditional probability* is the fraction of a row or column in a particular cell.

Example 6.2

Consider two random variables x and y , where x has five possible states and y has three possible states, as shown in Figure 6.2. The value c_i is the sum of the individual frequencies for the i^{th} column, that is $c_i = \sum_{j=1}^3 n_{ij}$. Similarly, the value r_j is the row sum, that is $r_j = \sum_{i=1}^5 n_{ij}$. Using these definitions, we can compactly express the distribution of x and y by themselves.

The probability distribution of each random variable, the marginal probability, which can be seen as the sum over a row or column

$$P(x = x_i) = \frac{c_i}{N} = \frac{\sum_{j=1}^3 n_{ij}}{N} \quad (6.10)$$

and

$$P(y = y_j) = \frac{r_j}{N} = \frac{\sum_{i=1}^5 n_{ij}}{N}, \quad (6.11)$$

where c_i and r_j are the i th column and j th row of the probability table, respectively. By convention for discrete random variables with a finite number of events, we assume that probabilities sum up to one, that is

$$\sum_{i=1}^5 P(x = x_i) = 1 \quad \text{and} \quad \sum_{j=1}^3 P(y = y_j) = 1. \quad (6.12)$$

The conditional probability is the fraction of a row or column in a particular cell. For example, the conditional probability of y given x is

$$P(y = y_j | x = x_i) = \frac{n_{ij}}{c_i}, \quad (6.13)$$

and the conditional probability of x given y is

$$P(x = x_i | y = y_j) = \frac{n_{ij}}{r_j}. \quad (6.14)$$

categorical variables

The marginal probability that x takes the value x irrespective of the value of random variable y is (lazily) written as $p(x)$. If we consider only the instances where $x = x$, then the fraction of instances (the conditional probability) for which $y = y$ is written (lazily) as $p(y | x)$.

In machine learning, we use discrete probability distributions to model *categorical variables*, i.e., variables that take a finite set of unordered values. These could be categorical features such as the degree taken at university when used for predicting the salary of a person, or categorical la-

bels such as letters of the alphabet when doing handwritten recognition. Discrete distributions are also often used to construct probabilistic models that combine a finite number of continuous distributions. We will see the Gaussian mixture model in Chapter 11.

6.2.2 Continuous Probabilities

We consider real valued random variables in this section, that is we consider state spaces which are intervals of the real line \mathbb{R} . In this book we will pretend that we can perform operations on real random variables as if we have discrete probability spaces with finite states. However this simplification is not precise for two situations: when we repeat something infinitely often, and when we want to draw a point from an interval. The first situation arises when we discuss generalization error in machine learning (Chapter 8). The second situation arises when we want to discuss continuous distributions such as the Gaussian (Section 6.5). For our purposes, the lack of precision allows a more brief introduction to probability.

Remark. In continuous spaces there are two additional technicalities which are counterintuitive. First the set of all subsets (used to define the event space \mathcal{A} in Section 6.1) is not well behaved enough. \mathcal{A} needs to be restricted to behave well under set complements, set intersections and set unions. Second the size of a set (which in discrete spaces can be obtained by counting the elements) turns out to be tricky. The size of a set is called its measure, for example the cardinality of discrete sets, the length of an interval in \mathbb{R} and the volume of a region in \mathbb{R}^d are all measures. Sets that behave well under set operations and furthermore have a topology are called a Borel σ -algebras. Betancourt (2018) details a careful construction of probability spaces from set theory without being bogged down in technicalities. A reader interested in a more precise construction is referred to Jacod and Protter (2004); Billingsley (1995). Further references can be found in the further reading section. In this book, we consider real valued random variables with their corresponding Borel σ -algebra. We consider random variables with values in \mathbb{R}^d to be a vector of real valued random variables. \diamond

Definition 6.1 (Probability Density Function). A function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is called a *probability density function* (pdf) if

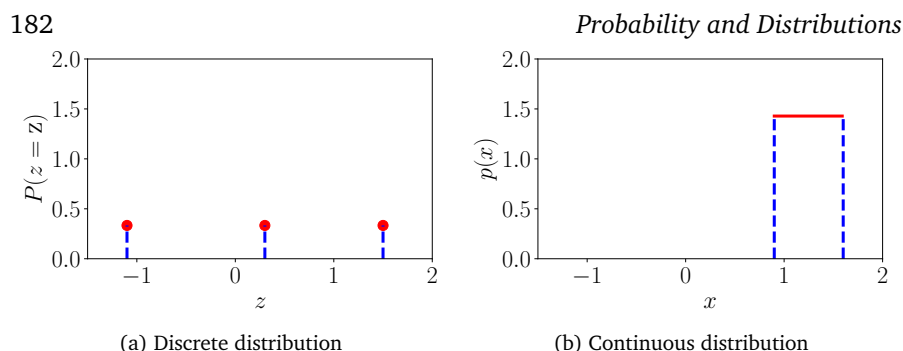
probability density
function

- 1 $\forall \mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) \geq 0$
- 2 Its integral exists and

$$\int_{\mathbb{R}^D} f(\mathbf{x}) d\mathbf{x} = 1. \quad (6.15)$$

Here, $\mathbf{x} \in \mathbb{R}^D$ is a (continuous) random variable. For probability mass functions (pmf) of discrete random variables the integral in (6.15) is replaced with a sum (see (6.12)).

Figure 6.3
Examples of discrete
and continuous
Uniform
distributions. See
Example 6.3 for
details of the
distributions.



$P(x = x)$ is a set of
measure zero.)

In contrast to discrete random variables, the probability of a continuous random variable x taking a particular value $P(x = x)$ is zero. However this does not mean that such events never occur.

cumulative
distribution function

Definition 6.2 (Cumulative Distribution Function). A *cumulative distribution function* (cdf) of a multivariate real-valued random variable $x \in \mathbb{R}^D$ is given by

$$F_x(x) = P(x_1 \leq x_1, \dots, x_D \leq x_D), \quad (6.16)$$

where the right-hand side represents the probability that random variable x_i takes the value smaller than or equal to x_i . This can be expressed also as the integral of the probability density function so that

$$F_x(x) = \int_{-\infty}^x f(x) dx. \quad (6.17)$$

6.2.3 Contrasting Discrete and Continuous Distributions

Let us consider both discrete and continuous distributions and contrast them. The aim here is to see that while both discrete and continuous distributions seem to have similar requirements, such as the total probability mass is 1, they are subtly different. Since the total probability mass of a discrete random variable is 1 (see (6.12)) and there are a finite number of states, the probability of each state must be in the interval $[0, 1]$. However, the analogous requirement for continuous random variables (see (6.15)) does not imply that the value of the density is less than or equal to 1 for all values. We illustrate this in Figure 6.3 using the *uniform distribution* for both discrete and continuous random variables.

Example 6.3

We consider two examples of the uniform distribution, where each state is equally likely to occur. This example illustrates the difference between discrete and continuous probability distributions.

Let z be a discrete uniform random variable with three states $\{z = -1.1, z = 0.3, z = 1.5\}$. Note that the actual values of these states are not meaningful here, and we deliberately chose numbers to drive home the point that we do not want to use (and should ignore) the ordering of the states. The probability mass function can be represented as a table of probability values.

z	-1.1	0.3	1.5
$P(z = Z)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

Alternatively, one could think of this as a graph (Figure 6.3(a)), where we use the fact that the states can be located on the x -axis, and the y -axis represents the probability of a particular state. The y -axis in Figure 6.3(a) is deliberately extended so that it is the same as in Figure 6.3(b).

Let x be a continuous random variable taking values in the range $0.9 \leq x \leq 1.6$, as represented by the graph in Figure 6.3(b). Observe that the height of the density can be greater than 1. However, it needs to hold that

$$\int_{0.9}^{1.6} p(x) dx = 1. \quad (6.18)$$

Remark. There is an additional subtlety with regards to discrete probability distributions. The states x_1, \dots, x_d do not in principle have any structure, that is there is usually no way to compare them, for example $x_1 = \text{red}, x_2 = \text{green}, x_3 = \text{blue}$. However there are applications where the discrete states form an ordered set, for example $x_1 = -1.1, x_2 = 0.3, x_3 = 1.5$, where we could say $x_1 < x_2 < x_3$. \diamond

Very often the literature uses notation and nomenclature that can be confusing to a beginner. For a value x of a state space Ω , $p(x)$ denotes the probability that random variable x takes value x , i.e., $P(x = x)$, which is known as the probability mass function. This is often referred to as the “distribution”. For continuous variables, $p(x)$ is called the probability density function (often referred to as a density), and to muddy things even further the cumulative distribution function $P(x \leq x)$ is often also referred to as the “distribution”. In this chapter, we often will use the notation x or \mathbf{x} to refer to univariate and multivariate random variables respectively. We summarise the nomenclature in Table 6.1.

Remark. We will be using the expression “probability distribution” not only for discrete probability mass functions but also for continuous probability density functions, although this is technically incorrect. Unfortunately the majority of machine learning literature is also sloppy about the phrase. \diamond

Table 6.1
Nomenclature for
probability
distributions.

	“point probability”	“interval probability”
discrete	$P(x = X)$ probability mass function	not applicable
continuous	$p(x)$ probability density function	$P(x \leq X)$ cumulative distribution function

6.3 Sum Rule, Product Rule and Bayes’ Theorem

When we think of a probabilistic model as an extension to logical reasoning, as we discussed in Section 6.1.1, the rules of probability presented here follow naturally from fulfilling the desiderata (Jaynes, 2003, Chapter 2). Probabilistic modeling provides a principled foundation for designing machine learning methods. Once we have defined probability distributions (Section 6.2) corresponding to the uncertainties of the data and our problem, it turns out that there are only two fundamental rules, the sum rule and the product rule, that govern probabilistic inference.

Given the definitions of marginal and conditional probability for discrete and continuous random variables in the previous section, we can now present the two fundamental rules in probability theory. These two rules arise naturally (Jaynes, 2003) from the requirements we discussed in Section 6.1.1. Recall from (6.9) that $p(\mathbf{x}, \mathbf{y})$ is the joint distribution of the two random variables \mathbf{x}, \mathbf{y} , $p(\mathbf{x})$, $p(\mathbf{y})$ are the corresponding marginal distributions, and $p(\mathbf{y} | \mathbf{x})$ is the conditional distribution of \mathbf{y} given \mathbf{x} .

The first rule, the *sum rule*, is

$$p(\mathbf{x}) = \begin{cases} \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{x}, \mathbf{y}) & \text{if } \mathbf{y} \text{ is discrete} \\ \int_{\mathcal{Y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} & \text{if } \mathbf{y} \text{ is continuous} \end{cases} \quad (6.19)$$

so that we sum out (or integrate out) the set of states \mathcal{Y} of the random variable \mathbf{y} . The sum rule is also known as the *marginalization property*. The sum rule relates the joint distribution to a marginal distribution. In general, when the joint distribution contains more than two random variables, the sum rule can be applied to any subset of the random variables, resulting in a marginal distribution of potentially more than one random variable. More concretely, if $\mathbf{x} = (x_1, \dots, x_D)$, we obtain the marginal

$$p(x_i) = \int p(x_1, \dots, x_D) d\mathbf{x}_{\setminus i} \quad (6.20)$$

by repeated application of the sum rule where we integrate/sum out all random variables except x_i , which is indicated by $\setminus i$.

Remark. Many of the computational challenges of probabilistic modeling are due to the application of the sum rule. When there are many variables or discrete variables with many states, the sum rule boils down to performing a high-dimensional sum or integral. Performing high dimensional

sums or integrals is generally computationally hard, in the sense that there is no known polynomial time algorithm to calculate them exactly. \diamond

The second rule, known as the *product rule*, relates the joint distribution to the conditional distribution via

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}). \quad (6.21)$$

The product rule can be interpreted as the fact that every joint distribution of two random variables can be factorized (written as a product) of two other distributions. The two factors are the marginal distribution of the first random variable $p(\mathbf{x})$, and the conditional distribution of the second random variable given the first $p(\mathbf{y} | \mathbf{x})$. Since the ordering of random variables is arbitrary in $p(\mathbf{x}, \mathbf{y})$ the product rule also implies $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x} | \mathbf{y})p(\mathbf{y})$. To be precise, (6.21) is expressed in terms of the probability mass functions for discrete random variables. For continuous random variables, the product rule is expressed in terms of the probability density functions (Section 6.2.3).

Let us briefly explore how to use probabilistic models to capture uncertainty (Ghahramani, 2015). At the lowest modeling level, measurement noise introduces model uncertainty, for example a camera sensor introduces random error in the value of each pixel it records. We will see in Chapter 9 how to use Gaussian (Section 6.5) noise models for linear regression. At higher modeling levels, we would be interested in modeling the uncertainty of the coefficients in linear regression. This uncertainty captures which values of these parameters will be good at predicting new data. Finally at the highest levels, we may want to capture uncertainties about the model structure. We will discuss model choice in Section 8.5. Once we have the probabilistic models (described in Section 8.3), the basic rules of probability presented in this section are used to infer the unobserved quantities given the observed data.

In machine learning and Bayesian statistics, we are often interested in making inferences of unobserved (latent) random variables given that we have observed other random variables. Let us assume we have some prior knowledge $p(\mathbf{x})$ about an unobserved random variable \mathbf{x} and some relationship $p(\mathbf{y} | \mathbf{x})$ between \mathbf{x} and a second random variable \mathbf{y} , which we can observe. If we observe \mathbf{y} we can use Bayes' theorem to draw some conclusions about \mathbf{x} given the observed values of \mathbf{y} . *Bayes' theorem* (also: *Bayes' rule* or *Bayes' law*)

Bayes' theorem is also called the "probabilistic inverse" Bayes' theorem

$$\underbrace{p(\mathbf{x} | \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y} | \mathbf{x})}^{\text{likelihood}} \overbrace{p(\mathbf{x})}^{\text{prior}}}{\underbrace{p(\mathbf{y})}_{\text{evidence}}} \quad (6.22)$$

is a direct consequence of the product rule in (6.19) since

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x} | \mathbf{y})p(\mathbf{y}) \quad (6.23)$$

and

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}) \quad (6.24)$$

so that

$$p(\mathbf{x} | \mathbf{y})p(\mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}) \iff p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}. \quad (6.25)$$

prior 3334 In (6.22), $p(\mathbf{x})$ is the *prior*, which encapsulates our subjective prior
 3335 knowledge of the unobserved (latent) variable \mathbf{x} before observing any
 3336 data. We can choose any prior that makes sense to us, but it is critical to
 3337 ensure that the prior has a non-zero pdf (or pmf) on all plausible \mathbf{x} , even
 3338 if they are very rare.

likelihood 3339 The *likelihood* $p(\mathbf{y} | \mathbf{x})$ describes how \mathbf{x} and \mathbf{y} are related, and it is the
 The likelihood is 3340 probability of the data \mathbf{y} if we were to know the latent variable \mathbf{x} . Note
 sometimes also 3341 that the likelihood is not a distribution in \mathbf{x} , but only in \mathbf{y} . We call $p(\mathbf{y} | \mathbf{x})$
 called the 3342 either the “likelihood of \mathbf{x} (given \mathbf{y})” or the “probability of \mathbf{y} given \mathbf{x} ” but
 “measurement 3343 never the likelihood of \mathbf{y} (MacKay, 2003a).
 model”.
 posterior 3344 The *posterior* $p(\mathbf{x} | \mathbf{y})$ is the quantity of interest in Bayesian statistics
 3345 because it expresses exactly what we are interested in, i.e., what we know
 3346 about \mathbf{x} after having observed \mathbf{y} .

The quantity

$$p(\mathbf{y}) := \int p(\mathbf{y} | \mathbf{x})p(\mathbf{x})d\mathbf{x} = \mathbb{E}_{\mathbf{x}}[p(\mathbf{y} | \mathbf{x})] \quad (6.26)$$

marginal likelihood 3347 is the *marginal likelihood/evidence*. By definition the marginal likelihood
 evidence 3348 integrates the numerator of (6.22) with respect to the latent variable \mathbf{x} .
 3349 Therefore, the marginal likelihood is independent of \mathbf{x} and it ensures
 3350 that the posterior $p(\mathbf{x} | \mathbf{y})$ is normalized. The marginal likelihood can also
 3351 be interpreted as the expected likelihood where we take the expectation
 3352 with respect to the prior $p(\mathbf{x})$. Beyond normalization of the posterior the
 3353 marginal likelihood also plays an important role in Bayesian model selec-
 3354 tion as we will discuss in Section 8.5. Due to the integration in (8.41), the
 3355 evidence is often hard to compute.

Bayes’ theorem in (6.22) allows us to invert the causal relationship be-
 tween \mathbf{x} and \mathbf{y} given by the likelihood. Therefore, Bayes’ theorem is some-
 times called the *probabilistic inverse*.

3359 *Remark.* In Bayesian statistics, the posterior distribution is the quantity
 3360 of interest as it encapsulates all available information from the prior and
 3361 the data. Instead of carrying the posterior around, it is possible to fo-
 3362 cus on some statistic of the posterior, such as the maximum of the pos-
 3363 terior, which we will discuss in Section 9.2.3. However, focusing on the
 3364 some statistic of the posterior leads to loss of information. If we think in
 3365 a bigger context, then the posterior can be used within a decision mak-
 3366 ing system, and having the full posterior around can be extremely useful
 3367 and lead to decisions that are robust to disturbances. For example, in the

context of model-based reinforcement learning, Deisenroth et al. (2015) show that using the full posterior distribution of plausible transition functions leads to very fast (data/sample efficient) learning, whereas focusing on the maximum of the posterior leads to consistent failures. Therefore, having the full posterior around in a downstream task can be very useful. In Chapter 9, we will continue this discussion in the context of linear regression. \diamond

6.4 Summary Statistics and Independence

We are often interested in summarizing sets of random variables and comparing pairs of random variables. A statistic of a random variable is a deterministic function of that random variable. The summary statistics of a distribution provides one useful view of how a random variable behaves, and as the name suggests, provides numbers that summarize and characterize the distribution. We describe the mean and the variance, two well-known summary statistics. Then we discuss two ways to compare a pair of random variables: first how to say that two random variables are independent, and second how to compute an inner product between them.

6.4.1 Means and Covariances

Mean and (co)variance are often useful to describe properties of probability distributions (expected values and spread). We will see in Section 6.6 that there is a useful family of distributions (called the exponential family), where the statistics of the random variable capture all possible information.

The main tool we use to compute statistics of a random variable is its expected value with respect to a particular function.

Definition 6.3 (Expected value). The *expected value* of a function $g : \mathbb{R} \rightarrow \mathbb{R}$ of a univariate continuous random variable $x \sim p(x)$ is given by

$$\mathbb{E}_x[g(x)] = \int g(x)p(x)dx. \quad (6.27)$$

Correspondingly the expected value of a function g of a discrete random variable $x \sim p(x)$ is given by

$$\mathbb{E}_x[g(x)] = \sum_{x \in \mathcal{A}} g(x)p(x) \quad (6.28)$$

where \mathcal{A} is the event space.

Remark. We consider multivariate random variables \mathbf{x} as a finite vector of

univariate random variables $[x_1, \dots, x_n]^\top$. For multivariate random variables, we define the expected value element wise

$$\mathbb{E}_{\mathbf{x}}[g(\mathbf{x})] = \begin{bmatrix} \mathbb{E}_{x_1}[g(x_1)] \\ \vdots \\ \mathbb{E}_{x_D}[g(x_D)] \end{bmatrix} \in \mathbb{R}^D, \quad (6.29)$$

where the subscript \mathbb{E}_{x_d} indicates that we are taking the expected value with respect to the d^{th} element of the vector \mathbf{x} . \diamond

Definition 6.3 defines the meaning of the notation \mathbb{E}_x and $\mathbb{E}_{\mathbf{x}}$ as the operator indicating that we should take the integral with respect to the probability density (for continuous distributions) or the sum over all states (for discrete distributions). The definition of the mean (Definition 6.4), is a special case of the expected value, obtained by choosing g to be the identity function.

Definition 6.4 (Mean). The *mean* of a random variable $x \in \mathbb{R}^D$ is an average and defined as

$$\mathbb{E}_{\mathbf{x}}[\mathbf{x}] = \begin{bmatrix} \mathbb{E}_{x_1}[x_1] \\ \vdots \\ \mathbb{E}_{x_D}[x_D] \end{bmatrix} \in \mathbb{R}^D, \quad (6.30)$$

where

$$\mathbb{E}_{x_d}[x_d] := \begin{cases} \int x_d p(x_d) dx_d & \text{if } x_d \text{ has a continuous domain} \\ \sum_i p(x_d = i) & \text{if } x_d \text{ has a discrete domain} \end{cases} \quad (6.31)$$

for $d = 1, \dots, D$, where the subscript indicates the corresponding dimension of \mathbf{x} .

In one dimension, there are two other intuitive notions of “average”, which are the *median* and the *mode*. The median is the “middle” value if we sort the values, i.e., 50% of the values are greater than the median and 50% are smaller than the median. This idea can be generalized to continuous values by considering the value where the CDF (Definition 6.2) is 0.5. For distributions which are asymmetric or have long tails, the median provides an estimate of a typical value that is closer to human intuition than the mean value. Furthermore the median is more robust to outliers than the mean. The generalization of the median to higher dimensions is non-trivial as there is no obvious way to “sort” in more than one dimension (Hallin et al., 2010; Kong and Mizera, 2012). The *mode* is the most frequently occurring value. For a discrete random variable, the mode is defined as the value of x having the highest frequency of occurrence. For a continuous random variable, the mode is defined as a peak in the density $p(\mathbf{x})$. A particular density $p(\mathbf{x})$ may have more than one mode,

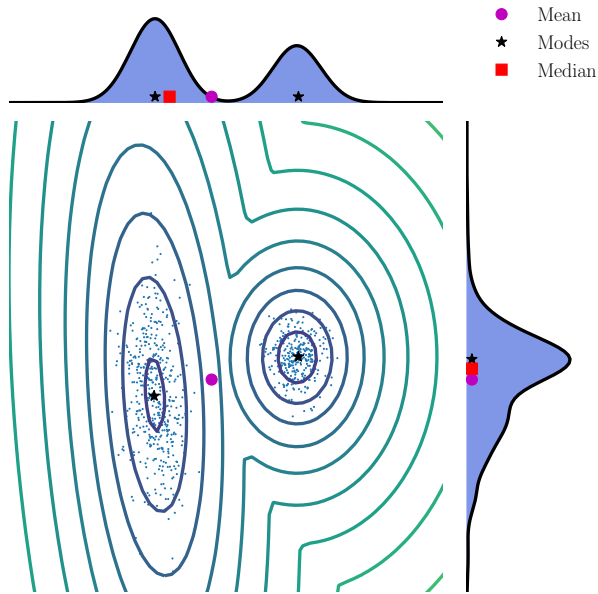


Figure 6.4
Illustration of the mean, mode and median for a two-dimensional dataset, as well as its marginal densities.

and, therefore, finding the mode may be computationally challenging in high dimensions.

Example 6.4

Consider the 2 dimensional distribution illustrated in Figure 6.4

$$\mathcal{N}\left(\begin{bmatrix} 10 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) + 1.5\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 & 2.9 \\ -1.1 & 0.7 \end{bmatrix}\right). \quad (6.32)$$

Also shown is its corresponding marginal distribution in each dimension. Observe that the distribution is bimodal (has two modes), but one of the marginal distributions is unimodal (has one mode). The horizontal bimodal univariate distribution illustrates the fact that the mean and median can be quite different from each other. While it is tempting to define the two dimensional median to be the concatenation of the medians in each dimension, the fact that we cannot define an ordering of two dimensional points makes it difficult. When we say cannot define an ordering, we mean that we there is more than one way to define $<$ such that $\begin{bmatrix} 1 \\ 0 \end{bmatrix} < \begin{bmatrix} 2 \\ 3 \end{bmatrix}$.

The mean is recovered if we set the function g in Definition 6.3 to the identity function. This indicates that we can think about functions of random variables, which we will revisit in Section 6.7.

Remark. The expected value is a linear operator. For example, given mul-

tivariate real-valued functions $f(\mathbf{x}) = ag(\mathbf{x}) + bh(\mathbf{x})$ where $a, b \in \mathbb{R}$,

$$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \quad (6.33a)$$

$$= \int [ag(\mathbf{x}) + bh(\mathbf{x})]p(\mathbf{x})d\mathbf{x} \quad (6.33b)$$

$$= a \int g(\mathbf{x})p(\mathbf{x})d\mathbf{x} + b \int h(\mathbf{x})p(\mathbf{x})d\mathbf{x} \quad (6.33c)$$

$$= a\mathbb{E}_{\mathbf{x}}[g(\mathbf{x})] + b\mathbb{E}_{\mathbf{x}}[h(\mathbf{x})]. \quad (6.33d)$$

◇

3424

3425

3426

3427

For two random variables, we may wish to characterize their correspondence to each other. The covariance intuitively represents the notion of how dependent random variables are to one another.

Definition 6.5 (Covariance (univariate)). The covariance between two univariate random variables $x, y \in \mathbb{R}$ is given by the expected product of their deviations from their respective means, that is

$$\text{Cov}[x, y] = \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])]. \quad (6.34)$$

By using the linearity of expectations, the expression in Definition 6.5 can be rewritten as the expected value of the product minus the product of the expected values, i.e.,

$$\text{Cov}[x, y] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]. \quad (6.35)$$

variance

3428

standard deviation

3429

3430

3431

3432

3433

3434

The covariance of a variable with itself $\text{Cov}[x, x]$ is called the *variance* and is denoted by $\mathbb{V}[x]$. The square root of the variance is called the *standard deviation* and is often denoted by $\sigma(x)$.

When we want to compare the covariances between different pairs of random variables, it turns out that the variance of each random variable affects the value of the covariance. The normalized version of covariance is called the correlation.

correlation

Definition 6.6 (Correlation). The *correlation* between two random variables x, y is given by

$$\text{corr}[x, y] = \frac{\text{Cov}[x, y]}{\sqrt{\mathbb{V}[x]\mathbb{V}[y]}}. \quad (6.36)$$

3435

3436

3437

3438

3439

3440

3441

3442

The correlation matrix is the covariance matrix of standardized random variables, $x/\sigma(x)$. In other words, each random variable is divided by its standard deviation (the square root of the variance) in the correlation matrix.

The covariance (and correlation) indicate how two random variables are related, see Figure 6.5. Positive correlation $\text{corr}[x, y]$ means that when x grows then y is also expected to grow. Negative correlation means that as x increases then y decreases.

Terminology:

Covariance of

multivariate random

variables $\text{Cov}[x, y]$

is sometimes

referred to as

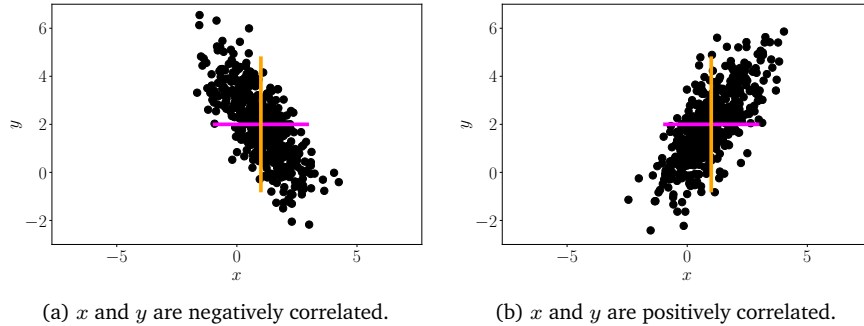
cross-covariance,

with covariance

referring to

 $\text{Cov}[x, x]$.

Figure 6.5
Two-dimensional
datasets with
identical means and
variances along
each axis (colored
lines) but with
different
correlations.



The notion of covariance can be generalised to multivariate random variables.

Definition 6.7 (Covariance). If we consider two random variables $\mathbf{x} \in \mathbb{R}^D$, $\mathbf{y} \in \mathbb{R}^E$, the *covariance* between \mathbf{x} and \mathbf{y} is defined as

covariance

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}_{\mathbf{x}}[\mathbf{x}]\mathbb{E}_{\mathbf{y}}[\mathbf{y}]^\top = \text{Cov}[\mathbf{y}, \mathbf{x}]^\top \in \mathbb{R}^{D \times E}. \quad (6.37)$$

Here, the subscript makes it explicit with respect to which variable we need to average.

Definition 6.7 can be applied with the same multivariate random variable in both arguments, which results in a useful concept that intuitively captures the “spread” of a random variable.

Definition 6.8 (Variance). The *variance* of a random variable $\mathbf{x} \in \mathbb{R}^D$ with mean vector $\boldsymbol{\mu} \in \mathbb{R}^D$ is defined as

variance

$$\mathbb{V}_{\mathbf{x}}[\mathbf{x}] = \mathbb{E}_{\mathbf{x}}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] = \mathbb{E}_{\mathbf{x}}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}_{\mathbf{x}}[\mathbf{x}]\mathbb{E}_{\mathbf{x}}[\mathbf{x}]^\top \quad (6.38a)$$

$$= \begin{bmatrix} \text{Cov}[x_1, x_1] & \text{Cov}[x_1, x_2] & \dots & \text{Cov}[x_1, x_D] \\ \text{Cov}[x_2, x_1] & \text{Cov}[x_2, x_2] & \dots & \text{Cov}[x_2, x_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[x_D, x_1] & \dots & \dots & \text{Cov}[x_D, x_D] \end{bmatrix} \in \mathbb{R}^{D \times D}. \quad (6.38b)$$

This matrix is called the *covariance matrix* of the random variable \mathbf{x} . The covariance matrix is symmetric and positive definite and tells us something about the spread of the data.

covariance matrix

The covariance matrix contains the variances of the *marginals*

marginal

$$p(x_i) = \int p(x_1, \dots, x_D) dx_{\setminus i} \quad (6.39)$$

on its diagonal, where “ $\setminus i$ ” denotes “all variables but i ”. The off-diagonal entries are the *cross-covariance* terms $\text{Cov}[x_i, x_j]$ for $i, j = 1, \dots, D$, $i \neq j$. It generally holds that

cross-covariance

$$\mathbb{V}_{\mathbf{x}}[\mathbf{x}] = \text{Cov}_{\mathbf{x}}[\mathbf{x}, \mathbf{x}]. \quad (6.40)$$

We will revisit the idea of covariance again in Section 6.4.5.

6.4.2 Empirical Means and Covariances

The definitions in Section 6.4.1 are often also called the *population mean and covariance*, as it refers to the true statistics for the population. In machine learning we need to learn from empirical observations of data. Consider a random variable x . There are two conceptual steps to go from population statistics to the realization of empirical statistics. First we use the fact that we have a finite dataset (of size N) to construct an empirical statistic which is a function of a finite number of identical random variables, x_1, \dots, x_N . Second we observe the data, that is we look at the realization of each of the random variables x_1, \dots, x_N and apply the empirical statistic.

Specifically for the mean (Definition 6.4), given a particular dataset we can obtain an estimate of the mean, which is called the *empirical mean* or *sample mean*. The same holds for the empirical covariance.

Definition 6.9 (Empirical Mean and Covariance). The *empirical mean* vector is the arithmetic average of the observations for each variable, and it is defined as

$$\bar{x} := \frac{1}{N} \sum_{n=1}^N x_n, \quad (6.41)$$

where $x_n \in \mathbb{R}^D$.

Similar to the empirical mean, the *empirical covariance* matrix is a $D \times D$ matrix

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^\top. \quad (6.42)$$

In the above definition, we have expressed the statistic in terms of the multivariate random variables x_1, \dots, x_N . To compute the statistics for a particular dataset, we would use the realizations (observations) x_1, \dots, x_N an use (6.41) and (6.42). Empirical covariance matrices are symmetric, positive semi-definite (see Section 3.2.3).

Remark. The notation we use in this book is imprecise in the sense that we do not explicitly make a distinction between the random variable x_n and its deterministic realization x_n . \diamond

6.4.3 Three Expressions for the Variance

We now focus on a single random variable x and use the empirical formulas above to derive three possible expressions for the variance. The derivation below is the same for the population variance, except that we need to take care of integrals. The standard definition of variance, corresponding

to the definition of covariance (Definition 6.5), is the expectation of the squared deviation of a random variable x from its expected value μ , i.e.,

$$\mathbb{V}_x[x] := \mathbb{E}_x[(x - \mu)^2]. \quad (6.43)$$

Depending on whether x is a discrete or continuous random variable, the expectation in (6.43) and the mean $\mu = \mathbb{E}_x(x)$ are computed using (6.31). The variance as expressed in (6.43) is the mean of a new random variable $z := (x - \mu)^2$.

When estimating the variance in (6.43) empirically, we need to resort to a two-pass algorithm: one pass through the data to calculate the mean μ using (6.41), and then a second pass using this estimate $\hat{\mu}$ calculate the variance. It turns out that we can avoid two passes by rearranging the terms. The formula in (6.43) can be converted to the so-called *raw-score formula for variance*:

raw-score formula
for variance

$$\mathbb{V}_x[x] = \mathbb{E}_x[x^2] - \mathbb{E}_x[x]^2. \quad (6.44)$$

The expression in (6.44) can be remembered as “the mean of the square minus the square of the mean”. It can be calculated empirically in one pass through data since we can accumulate x_i (to calculate the mean) and x_i^2 simultaneously. Unfortunately, if implemented in this way, it can be numerically unstable. The raw-score version of the variance can be useful in machine learning, e.g., when deriving the bias-variance decomposition (Bishop, 2006).

The two terms can
cancel out, resulting
in a loss of
numerical precision
in floating point
arithmetic.

A third way to understand the variance is that it is a sum of pairwise differences between all pairs of observations. Consider a sample x_1, \dots, x_N of realizations of random variable x , and we compute the squared difference between pairs of x_i and x_j . By expanding the square we can show that the sum of N^2 pairwise differences is the empirical variance of the observations,

$$\frac{1}{N^2} \sum_{i,j=1}^N (x_i - x_j)^2 = 2 \left[\frac{1}{N} \sum_{i=1}^N x_i^2 - \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2 \right] \quad (6.45)$$

We see that (6.45) is twice the raw-score expression (6.44). This means that we can express the sum of pairwise distances (of which there are N^2 of them) as a sum of deviations from the mean (of which there are N). Geometrically, this means that there is an equivalence between the pairwise distances and the distances from the center of the set of points. From a computational perspective, this means that by computing the mean (N terms in the summation), and then computing the variance (again N terms in the summation) we can obtain an expression (left-hand side of (6.45)) that has N^2 terms.

6.4.4 Sums and Transformations of Random Variables

We may want to model a phenomenon that cannot be well explained by textbook distributions (we introduce some in Sections 6.5 and 6.6), and hence may perform simple manipulations of random variables (such as adding two random variables).

Consider two random variables $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$. It holds that

$$\mathbb{E}[\mathbf{x} + \mathbf{y}] = \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{y}] \quad (6.46)$$

$$\mathbb{E}[\mathbf{x} - \mathbf{y}] = \mathbb{E}[\mathbf{x}] - \mathbb{E}[\mathbf{y}] \quad (6.47)$$

$$\mathbb{V}[\mathbf{x} + \mathbf{y}] = \mathbb{V}[\mathbf{x}] + \mathbb{V}[\mathbf{y}] + \text{Cov}[\mathbf{x}, \mathbf{y}] + \text{Cov}[\mathbf{y}, \mathbf{x}] \quad (6.48)$$

$$\mathbb{V}[\mathbf{x} - \mathbf{y}] = \mathbb{V}[\mathbf{x}] + \mathbb{V}[\mathbf{y}] - \text{Cov}[\mathbf{x}, \mathbf{y}] - \text{Cov}[\mathbf{y}, \mathbf{x}] \quad (6.49)$$

Mean and (co)variance exhibit some useful properties when it comes to affine transformation of random variables. Consider a random variable \mathbf{x} with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ and a (deterministic) affine transformation $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ of \mathbf{x} . Then \mathbf{y} is itself a random variable whose mean vector and covariance matrix are given by

$$\mathbb{E}_{\mathbf{y}}[\mathbf{y}] = \mathbb{E}_{\mathbf{x}}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{E}_{\mathbf{x}}[\mathbf{x}] + \mathbf{b} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \quad (6.50)$$

$$\mathbb{V}_{\mathbf{y}}[\mathbf{y}] = \mathbb{V}_{\mathbf{x}}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbb{V}_{\mathbf{x}}[\mathbf{A}\mathbf{x}] = \mathbf{A}\mathbb{V}_{\mathbf{x}}[\mathbf{x}]\mathbf{A}^{\top} = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^{\top}, \quad (6.51)$$

This can be shown directly by using the definition of the mean and covariance.

respectively. Furthermore,

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}[\mathbf{x}(\mathbf{A}\mathbf{x} + \mathbf{b})^{\top}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{A}\mathbf{x} + \mathbf{b}]^{\top} \quad (6.52)$$

$$= \mathbb{E}[\mathbf{x}]\mathbf{b}^{\top} + \mathbb{E}[\mathbf{x}\mathbf{x}^{\top}]\mathbf{A}^{\top} - \boldsymbol{\mu}\mathbf{b}^{\top} - \boldsymbol{\mu}\boldsymbol{\mu}^{\top}\mathbf{A}^{\top} \quad (6.53)$$

$$= \boldsymbol{\mu}\mathbf{b}^{\top} - \boldsymbol{\mu}\mathbf{b}^{\top} + (\mathbb{E}[\mathbf{x}\mathbf{x}^{\top}] - \boldsymbol{\mu}\boldsymbol{\mu}^{\top})\mathbf{A}^{\top} \quad (6.54)$$

$$\stackrel{(6.38a)}{=} \boldsymbol{\Sigma}\mathbf{A}^{\top}, \quad (6.55)$$

where $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top}] - \boldsymbol{\mu}\boldsymbol{\mu}^{\top}$ is the covariance of \mathbf{x} .

6.4.5 Statistical Independence

statistically independent

Definition 6.10 (Independence). Two random variables \mathbf{x}, \mathbf{y} are *statistically independent* if and only if

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}). \quad (6.56)$$

Intuitively, two random variables \mathbf{x} and \mathbf{y} are independent if the value of \mathbf{y} (once known) does not add any additional information about \mathbf{x} (and vice versa).

If \mathbf{x}, \mathbf{y} are (statistically) independent then

- $p(\mathbf{y} | \mathbf{x}) = p(\mathbf{y})$
- $p(\mathbf{x} | \mathbf{y}) = p(\mathbf{x})$
- $\mathbb{V}_{\mathbf{x}, \mathbf{y}}[\mathbf{x} + \mathbf{y}] = \mathbb{V}_{\mathbf{x}}[\mathbf{x}] + \mathbb{V}_{\mathbf{y}}[\mathbf{y}]$
- $\text{Cov}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}, \mathbf{y}] = \mathbf{0}$

The last point above may not hold in converse, that is two random variables can have covariance zero but are not statistically independent. To understand why, recall that covariance measures only linear dependence, therefore random variables that are non-linearly dependent could have covariance zero.

Example 6.5

Consider a random variable x with zero mean ($\mathbb{E}_x[x] = 0$) and also $\mathbb{E}_x[x^3] = 0$. Let $y = x^2$ (hence y is dependent on x) and consider the covariance (6.35) between x and y . But this gives

$$\text{Cov}[x, y] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y] = \mathbb{E}[x^3] = 0. \quad (6.57)$$

In machine learning we often consider problems that can be modelled as *independent and identically distributed* random variables, x_1, \dots, x_N . The word “independent” refers to Definition 6.10, that is any pair of random variables x_i and x_j are independent. The phrase identically distributed means that all the random variables are from the same distribution.

independent and
identically
distributed

Another concept that is important in machine learning is conditional independence.

Definition 6.11 (Conditional Independence). Formally, two random variables \mathbf{x} and \mathbf{y} are *conditionally independent given \mathbf{z}* if and only if

$$p(\mathbf{x}, \mathbf{y} | \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{y} | \mathbf{z}) \quad \text{for all } \mathbf{z} \in \mathcal{A}. \quad (6.58)$$

conditionally
independent given \mathbf{z}

We write $\mathbf{x} \perp\!\!\!\perp \mathbf{y} | \mathbf{z}$.

Definition 6.11 requires that the relation in (6.58) must hold true for every value of \mathbf{z} . The interpretation of (6.58) can be understood as “given knowledge about \mathbf{z} , the distribution of \mathbf{x} and \mathbf{y} factorizes”. Independence can be cast as a special case of conditional independence if we write $\mathbf{x} \perp\!\!\!\perp \mathbf{y} | \emptyset$.

By using the product rule of probability from (6.21) we can expand the left-hand side of (6.58) to obtain

$$p(\mathbf{x}, \mathbf{y} | \mathbf{z}) = p(\mathbf{x} | \mathbf{y}, \mathbf{z})p(\mathbf{y} | \mathbf{z}). \quad (6.59)$$

By comparing the right-hand side of (6.58) with (6.59) we see that $p(\mathbf{y} | \mathbf{z})$ appears in both of them so that

$$p(\mathbf{x} | \mathbf{y}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z}). \quad (6.60)$$

Equation (6.60) provides an alternative definition of conditional independence, i.e., $\mathbf{x} \perp\!\!\!\perp \mathbf{y} | \mathbf{z}$. This alternative presentation provides the interpretation “given that we know \mathbf{z} , knowledge about \mathbf{y} does not change our knowledge of \mathbf{x} ”.

6.4.6 Inner Products of Random Variables

Recall the definition of inner products from Section 3.2. Another example for defining an inner product between unusual types are random variables or random vectors. If we have two uncorrelated random variables x, y then

$$\mathbb{V}[x + y] = \mathbb{V}[x] + \mathbb{V}[y] \quad (6.61)$$

Since variances are measured in squared units, this looks very much like the Pythagorean theorem for right triangles $c^2 = a^2 + b^2$.

In the following, we see whether we can find a geometric interpretation of the variance relation of uncorrelated random variables in (6.61). Random variables can be considered vectors in a vector space, and we can define inner products to obtain geometric properties of random variables (Eaton, 2007). If we define

$$\langle x, y \rangle := \text{Cov}[x, y] \quad (6.62)$$

for zero mean random variables x and y , we obtain an inner product. we see that the covariance is symmetric, positive definite, and linear in either argument. The length of a random variable is

$$\|x\| = \sqrt{\text{Cov}[x, x]} = \sqrt{\mathbb{V}[x]} = \sigma[x], \quad (6.63)$$

i.e., its standard deviation. The “longer” the random variable, the more uncertain it is; and a random variable with length 0 is deterministic.

If we look at the angle θ between random two random variables x, y , we get

$$\cos \theta = \frac{\langle x, y \rangle}{\|x\| \|y\|} = \frac{\text{Cov}[x, y]}{\sqrt{\mathbb{V}[x]\mathbb{V}[y]}}, \quad (6.64)$$

which is the correlation (Definition 6.6) between the two random variables. This means that we can think of correlation as the angle between two random variables when we consider them geometrically. We know from Definition 3.7 that $x \perp y \iff \langle x, y \rangle = 0$. In our case this means that x and y are orthogonal if and only if $\text{Cov}[x, y] = 0$, i.e., they are uncorrelated. Figure 6.6 illustrates this relationship.

Remark. While it is tempting to use the Euclidean distance (constructed from the definition of inner products above) to compare probability distributions, it is unfortunately not the best way to obtain distances between distributions. Recall that the probability mass (or density) is positive and needs to add up to 1. These constraints mean that distributions live on something called a statistical manifold. The study of this space of probability distributions is called information geometry. Computing distances between distributions are often done using Kullback-Leibler divergence which is a generalization of distances that account for properties of the statistical manifold. Just like the Euclidean distance is a special case of a

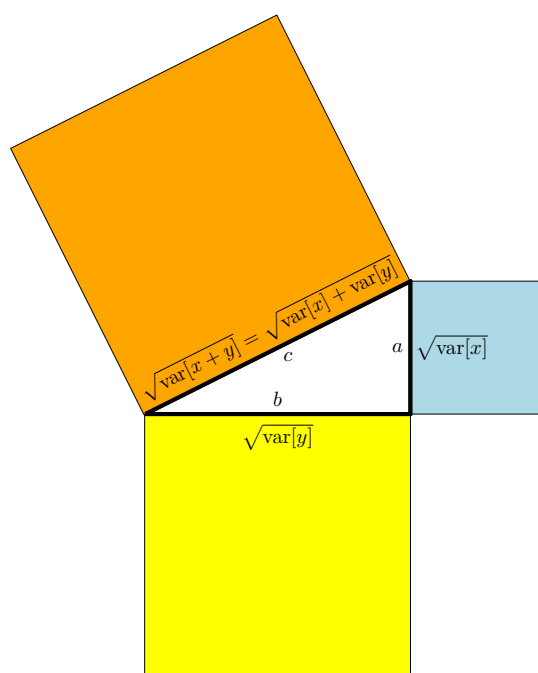


Figure 6.6
Geometry of random variables. If random variables x and y are uncorrelated they are orthogonal vectors in a corresponding vector space, and the Pythagorean theorem applies.

metric (Section 3.3) the Kullback-Leibler divergence is a special case of two more general classes of divergences called Bregman divergences and f -divergences. The study of divergences is beyond the scope of this book. Interested readers are referred to a recent book (Amari, 2016) written by one of the founders of the field of information geometry. \diamond

6.5 Gaussian Distribution

The Gaussian distribution is the most important probability distribution for continuous-valued random variables. It is also referred to as the *normal distribution*. Its importance originates from the fact that it has many computationally convenient properties, which we will be discussing in the following. In particular, we will use it to define the likelihood and prior for linear regression (Chapter 9), and consider a mixture of Gaussians for density estimation (Chapter 11).

There are many other areas of machine learning that also benefit from using a Gaussian distribution, for example Gaussian processes, variational inference and reinforcement learning. It is also widely used in other application areas such as signal processing (e.g., Kalman filter), control (e.g., linear quadratic regulator) and statistics (e.g. hypothesis testing).

For a univariate random variable, the Gaussian distribution has a den-

The Gaussian distribution arises naturally when we consider sums of independent and identically distributed random variables. This is known as the Central Limit Theorem (Grinstead and Snell, 1997). normal distribution

Figure 6.7
Gaussian
distribution of two
random variables
 x, y .

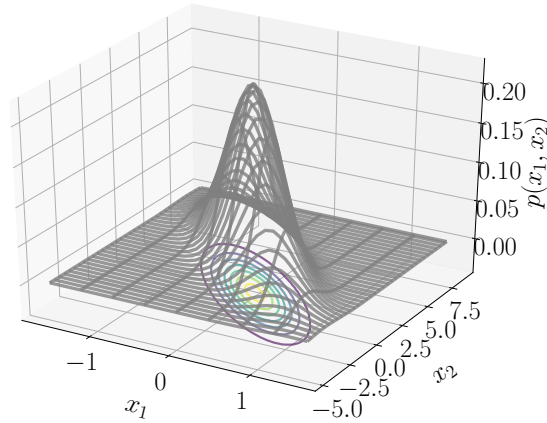
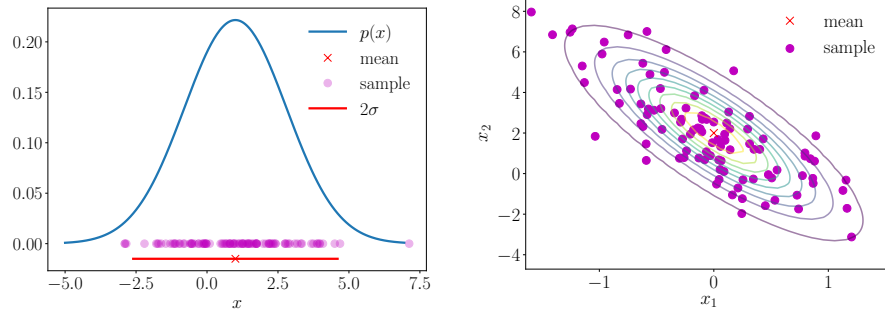


Figure 6.8
Gaussian
distributions
overlaid with 100
samples. Left:
Univariate
(1-dimensional)
Gaussian; The red
cross shows the
mean and the red
line shows the
extent of the
variance. Right:
Multivariate
(2-dimensional)
Gaussian, viewed
from top. The red
cross shows the
mean and the
coloured lines
shows the contour
lines of the density.



sity that is given by

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right). \quad (6.65)$$

The *multivariate Gaussian distribution* is fully characterized by a *mean vector* $\boldsymbol{\mu}$ and a *covariance matrix* $\boldsymbol{\Sigma}$ and defined as

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad (6.66)$$

where $\mathbf{x} \in \mathbb{R}^D$ is a random variable. We write $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ or $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Figure 6.7 shows a bi-variate Gaussian (mesh), with the corresponding contour plot. The special case of the Gaussian with zero mean and identity covariance, that is $\boldsymbol{\mu} = \mathbf{0}$ and $\boldsymbol{\Sigma} = \mathbf{I}$, is referred to as the *standard normal distribution*.

Gaussian distributions are widely used in statistical estimation and machine learning because they have closed-form expressions for marginal and conditional distributions. In Chapter 9, we use these closed form expressions extensively for linear regression. A major advantage of modelling with Gaussian distributed random variables is that variable trans-

formations (Section 6.7) are often not needed. Since the Gaussian distribution is fully specified by its mean and covariance we often can obtain the transformed distribution by applying the transformation to the mean and covariance of the random variable.

6.5.1 Marginals and Conditionals of Gaussians are Gaussians

In the following, we present marginalization and conditioning in the general case of multivariate random variables. If this is confusing at first reading, the reader is advised to consider two univariate random variables instead. Let \mathbf{x} and \mathbf{y} be two multivariate random variables, which may have different dimensions. We would like to consider the effect of applying the sum rule of probability and the effect of conditioning. We therefore explicitly write the Gaussian distribution in terms of the concatenated random variable $[\mathbf{x}, \mathbf{y}]^\top$,

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix}\right). \quad (6.67)$$

where $\boldsymbol{\Sigma}_{xx} = \text{Cov}[\mathbf{x}, \mathbf{x}]$ and $\boldsymbol{\Sigma}_{yy} = \text{Cov}[\mathbf{y}, \mathbf{y}]$ are the marginal covariance matrices of \mathbf{x} and \mathbf{y} , respectively, and $\boldsymbol{\Sigma}_{xy} = \text{Cov}[\mathbf{x}, \mathbf{y}]$ is the cross-covariance matrix between \mathbf{x} and \mathbf{y} .

The conditional distribution $p(\mathbf{x} | \mathbf{y})$ is also Gaussian (illustrated on the bottom right of Figure 6.9) and given by (derived in Section 2.3 of Bishop (2006))

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y}) \quad (6.68)$$

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) \quad (6.69)$$

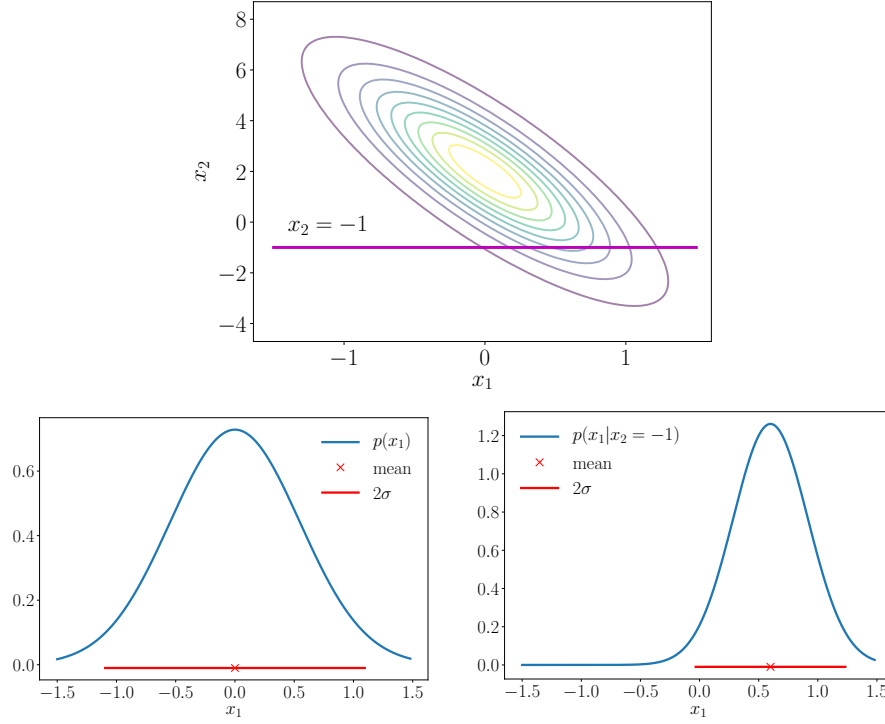
$$\boldsymbol{\Sigma}_{x|y} = \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx}. \quad (6.70)$$

Note that in the computation of the mean in (6.69) the \mathbf{y} -value is an observation and no longer random.

Remark. The conditional Gaussian distribution shows up in many places, where we are interested in posterior distributions:

- The Kalman filter (Kalman, 1960), one of the most central algorithms for state estimation in signal processing, does nothing but computing Gaussian conditionals of joint distributions (Deisenroth and Ohlsson, 2011).
- Gaussian processes (Rasmussen and Williams, 2006), which are a practical implementation of a distribution over functions. In a Gaussian process, we make assumptions of joint Gaussianity of random variables. By (Gaussian) conditioning on observed data, we can determine a posterior distribution over functions.
- Latent linear Gaussian models (Roweis and Ghahramani, 1999; Murphy, 2012), which include probabilistic PCA (Tipping and Bishop, 1999).

Figure 6.9 Top: Bivariate Gaussian; Bottom left: Marginal of a joint Gaussian distribution is Gaussian; Bottom right: The conditional distribution of a Gaussian is also Gaussian



3608

◇

The marginal distribution $p(\mathbf{x})$ of a joint Gaussian distribution $p(\mathbf{x}, \mathbf{y})$, see (6.67), is itself Gaussian and computed by applying the sum-rule in (6.19) and given by

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}). \quad (6.71)$$

3609

The corresponding result holds for $p(\mathbf{y})$, which is obtained by marginalizing with respect to \mathbf{x} . Intuitively, looking at the joint distribution in (6.67), we ignore (i.e., integrate out) everything we are not interested in. This is illustrated on the bottom left of Figure 6.9.

3610

3611

3612

Example 6.6

Consider the bivariate Gaussian distribution (illustrated in Figure 6.9)

$$p(x, y) = \mathcal{N}\left(\begin{bmatrix} 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 0.3 & -1 \\ -1 & 5 \end{bmatrix}\right). \quad (6.72)$$

We can compute the parameters of the univariate Gaussian, conditioned on $y = -1$, by applying (6.69) and (6.70) to obtain the mean and variance respectively. Numerically, this is

$$\mu_{x|y=-1} = 0 + (-1)(0.2)(-1 - 2) = 0.6 \quad (6.73)$$

and

$$\sigma_{x|y=-1}^2 = 0.3 - (-1)(0.2)(-1) = 0.1. \quad (6.74)$$

Therefore the conditional Gaussian is given by

$$p(x | y = -1) = \mathcal{N}(0.6, 0.1). \quad (6.75)$$

The marginal distribution $p(x)$ in contrast can be obtained by applying (6.71), which is essentially using the mean and variance of the random variable x , giving us

$$p(x) = \mathcal{N}(0, 0.3) \quad (6.76)$$

6.5.2 Product of Gaussian Densities

For linear regression (Chapter 9), we need to compute a Gaussian likelihood. Furthermore we may wish to assume a Gaussian prior (Section 9.3). The application of Bayes rule to compute the posterior results in a multiplication of the likelihood and the prior, that is the multiplication of two Gaussian densities. The *product* of two Gaussians $\mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A})\mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B})$ is a Gaussian distribution scaled by a $c \in \mathbb{R}$, given by $c\mathcal{N}(\mathbf{x} | \mathbf{c}, \mathbf{C})$ with

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \quad (6.77)$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b}) \quad (6.78)$$

$$c = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1} (\mathbf{a} - \mathbf{b})\right). \quad (6.79)$$

The scaling constant c itself can be written in the form of a Gaussian density either in \mathbf{a} or in \mathbf{b} with an “inflated” covariance matrix $\mathbf{A} + \mathbf{B}$, i.e., $c = \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B}) = \mathcal{N}(\mathbf{b} | \mathbf{a}, \mathbf{A} + \mathbf{B})$.

Remark. For notation convenience, we will sometimes use $\mathcal{N}(\mathbf{x} | \mathbf{m}, \mathbf{S})$ to describe the functional form of a Gaussian even if \mathbf{x} is not a random variable. We have just done this above when we wrote

$$c = \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B}) = \mathcal{N}(\mathbf{b} | \mathbf{a}, \mathbf{A} + \mathbf{B}). \quad (6.80)$$

Here, neither \mathbf{a} nor \mathbf{b} are random variables. However, writing c in this way is more compact than (6.79). \diamond

6.5.3 Sums and Linear Transformations

If \mathbf{x}, \mathbf{y} are independent Gaussian random variables (i.e., the joint is given as $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$) with $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$ and $p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$, then $\mathbf{x} + \mathbf{y}$ is also Gaussian distributed and given by

$$p(\mathbf{x} + \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_x + \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_x + \boldsymbol{\Sigma}_y). \quad (6.81)$$

Knowing that $p(\mathbf{x} + \mathbf{y})$ is Gaussian, the mean and covariance matrix can be determined immediately using the results from (6.46)–(6.49). This property will be important when we consider i.i.d. Gaussian noise acting on random variables as is the case for linear regression (Chapter 9).

Example 6.7

Since expectations are linear operations, we can obtain the weighted sum of independent Gaussian random variables

$$p(a\mathbf{x} + b\mathbf{y}) = \mathcal{N}(a\boldsymbol{\mu}_x + b\boldsymbol{\mu}_y, a\boldsymbol{\Sigma}_x + b\boldsymbol{\Sigma}_y). \quad (6.82)$$

Remark. A case which will be useful in Chapter 11 is the weighted sum of Gaussian densities. This is different from the weighted sum of Gaussian random variables. \diamond

In Theorem 6.12, the random variable x is from a density which a mixture of two densities $p_1(x)$ and $p_2(x)$, weighted by α . The theorem can be generalized to the multivariate random variable case, since linearity of expectations holds also for multivariate random variables. However the idea of a squared random variable needs to be replaced by $\mathbf{x}\mathbf{x}^\top$.

Theorem 6.12. Consider a weighted sum of two univariate Gaussian densities

$$p(x) = \alpha p_1(x) + (1 - \alpha)p_2(x) \quad (6.83)$$

where the scalar $0 < \alpha < 1$ is the mixture weight, and $p_1(x)$ and $p_2(x)$ are univariate Gaussian densities (Equation (6.65)) with different parameters, that is $(\mu_1, \sigma_1^2) \neq (\mu_2, \sigma_2^2)$.

The mean of the mixture x is given by the weighted sum of the means of each random variable,

$$\mathbb{E}[x] = \alpha\mu_1 + (1 - \alpha)\mu_2. \quad (6.84)$$

The variance of the mixture x is the mean of the conditional variance and the variance of the conditional mean,

$$\mathbb{V}[x] = [\alpha\sigma_1^2 + (1 - \alpha)\sigma_2^2] + \left([\alpha\mu_1^2 + (1 - \alpha)\mu_2^2] - [\alpha\mu_1 + (1 - \alpha)\mu_2]^2 \right). \quad (6.85)$$

Proof The mean of the mixture x is given by the weighted sum of the means of each random variable. We apply the definition of the mean (Definition 6.4), and plug in our mixture (Equation (6.83)) above

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} xp(x)dx \quad (6.86a)$$

$$= \int_{-\infty}^{\infty} \alpha xp_1(x) + (1 - \alpha)xp_2(x)dx \quad (6.86b)$$

$$= \alpha \int_{-\infty}^{\infty} xp_1(x)dx + (1 - \alpha) \int_{-\infty}^{\infty} xp_2(x)dx \quad (6.86c)$$

$$= \alpha\mu_1 + (1 - \alpha)\mu_2. \quad (6.86d)$$

To compute the variance, we can use the raw score version of the variance (Equation (6.44)), which requires an expression of the expectation of the squared random variable. Here we use the definition of an expectation of a function (the square) of a random variable (Definition 6.3).

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} x^2 p(x) dx \quad (6.87a)$$

$$= \int_{-\infty}^{\infty} \alpha x^2 p_1(x) + (1 - \alpha) x^2 p_2(x) dx \quad (6.87b)$$

$$= \alpha \int_{-\infty}^{\infty} x^2 p_1(x) dx + (1 - \alpha) \int_{-\infty}^{\infty} x^2 p_2(x) dx \quad (6.87c)$$

$$= \alpha(\mu_1^2 + \sigma_1^2) + (1 - \alpha)(\mu_2^2 + \sigma_2^2). \quad (6.87d)$$

where in the last equality, we again used the raw score version of the variance and rearranged terms such that the expectation of a squared random variable is the sum of the squared mean and the variance.

Therefore the variance is given by subtracting (6.86d) from (6.87d),

$$\mathbb{V}[x] = \mathbb{E}[x^2] - (\mathbb{E}[x])^2 \quad (6.88a)$$

$$= \alpha(\mu_1^2 + \sigma_1^2) + (1 - \alpha)(\mu_2^2 + \sigma_2^2) - (\alpha\mu_1 + (1 - \alpha)\mu_2)^2 \quad (6.88b)$$

$$= [\alpha\sigma_1^2 + (1 - \alpha)\sigma_2^2] + \left([\alpha\mu_1^2 + (1 - \alpha)\mu_2^2] - [\alpha\mu_1 + (1 - \alpha)\mu_2]^2 \right). \quad (6.88c)$$

For a mixture, the individual components can be considered to be conditional distributions (conditioned on the component identity). The last line is an illustration of the conditional variance formula: “The variance of a mixture is the mean of the conditional variance and the variance of the conditional mean”. \square

Remark. The derivation above holds for any density, but in the case of the Gaussian since it is fully determined by the mean and variance, the mixture density can be determined in closed form. \diamond

We consider in Example 6.18 a bivariate standard Gaussian random variable \mathbf{x} and performed a linear transformation $\mathbf{A}\mathbf{x}$ on it. The outcome is a Gaussian random variable with zero mean and covariance $\mathbf{A}\mathbf{A}^\top$. Observe that adding a constant vector will change the mean of the distribution, without affecting its variance, that is the random variable $\mathbf{x} + \boldsymbol{\mu}$ is Gaussian with mean $\boldsymbol{\mu}$ and identity covariance. Therefore, a linear (or affine) transformation of a Gaussian random variable is Gaussian distributed.

Consider a Gaussian distributed random variable $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$. For

a given matrix \mathbf{A} of appropriate shape, let \mathbf{y} be a random variable $\mathbf{y} = \mathbf{A}\mathbf{x}$ which is a transformed version of \mathbf{x} . We can compute the mean of \mathbf{y} by using the fact that the expectation is a linear operator (Equation (6.50)) as follows:

$$\mathbb{E}[\mathbf{A}\mathbf{x}] = \mathbf{A}\mathbb{E}[\mathbf{x}] = \mathbf{A}\boldsymbol{\mu}. \quad (6.89)$$

Similarly the variance of \mathbf{y} can be found by using Equation (6.51):

$$\mathbb{V}[\mathbf{A}\mathbf{x}] = \mathbf{A}\mathbb{V}[\mathbf{x}]\mathbf{A}^\top = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top. \quad (6.90)$$

This means that the random variable \mathbf{y} is distributed according to

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{x} | \mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top). \quad (6.91)$$

Let us now consider the reverse transformation: when we know that a random variable has a mean that is a linear transformation of another random variable. For a given full rank matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ where $m \geq n$, let $\mathbf{y} \in \mathbb{R}^m$ be a Gaussian random variable with mean $\mathbf{A}\mathbf{x}$, i.e.,

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x}, \boldsymbol{\Sigma}). \quad (6.92)$$

What is the corresponding probability distribution $p(\mathbf{x})$? If \mathbf{A} is invertible, then we can write $\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$ and apply the transformation in the previous paragraph. However, in general \mathbf{A} is not invertible, and we use an approach similar to that of the pseudo-inverse (3.56). That is we pre-multiply both sides with \mathbf{A}^\top and then invert $\mathbf{A}^\top\mathbf{A}$ which is symmetric and positive definite, giving us the relation

$$\mathbf{y} = \mathbf{A}\mathbf{x} \iff (\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top\mathbf{y} = \mathbf{x}. \quad (6.93)$$

Hence, \mathbf{x} is a linear transformation of \mathbf{y} , and we obtain

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | (\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top\mathbf{y}, (\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top\boldsymbol{\Sigma}\mathbf{A}(\mathbf{A}^\top\mathbf{A})^{-1}). \quad (6.94)$$

6.5.4 Sampling from Multivariate Gaussian Distributions

We will not explain the subtleties of random sampling on a computer. In the case of a multivariate Gaussian, this process consists of three stages: first we need a source of pseudo-random numbers that provide a uniform sample in the interval $[0,1]$, second we use a non-linear transformation such as the Box-Müller transform (Devroye, 1986) to obtain a sample from a univariate Gaussian, and third we collate a vector of these samples to obtain a sample from a multivariate standard normal $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

For a general multivariate Gaussian, that is where the mean is non-zero and the covariance is not the identity matrix, we use the properties of linear transformations of a Gaussian random variable. Assume we are interested in generating samples $\mathbf{x}_i, i = 1, \dots, n$, from a multivariate Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. We would like to construct the sample from a sampler that provides samples from the multivariate standard normal $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

To compute the Cholesky factorization of a matrix, it is required that the matrix is symmetric and positive definite (Section 3.2.3). Covariance matrices possess this property.

Draft (2018-11-13) from Mathematics for Machine Learning. Errata and feedback to <https://mml-book.com>.

To obtain samples from a multivariate normal $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we can use the properties of a linear transformation of a Gaussian random variable: If $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ then $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\mu}$, where $\mathbf{A}\mathbf{A}^\top = \boldsymbol{\Sigma}$, is Gaussian distributed with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Recall from Section 4.3 that when we can decompose $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^\top$, while there are many possible decompositions, we often choose the Cholesky decomposition. This has the benefit that \mathbf{A} is triangular, leading to efficient computation.

6.6 Conjugacy and the Exponential Family

Many of the probability distributions “with names” that we find in statistics textbooks were discovered to model particular types of phenomena. For example we have seen the Gaussian distribution in Section 6.5. The distributions are also related to each other in complex ways (Leemis and McQueston, 2008). For a beginner in the field, it can be overwhelming to figure out which distribution to use. In addition, many of these distributions were discovered at a time that statistics and computation was done by pencil and paper. It is natural to ask what are meaningful concepts in the computing age (Efron and Hastie, 2016). In the previous section, we saw that many of the operations required for inference can be conveniently calculated when the distribution is Gaussian. It is worth recalling at this point the desiderata for manipulating probability distributions.

- 1 There is some “closure property” when applying the rules of probability, e.g., Bayes’ theorem. By closure we mean that applying a particular operation returns an object of the same type.
- 2 As we collect more data, we do not need more parameters to describe the distribution.
- 3 Since we are interested in learning from data, we want parameter estimation to behave nicely.

It turns out that the class of distributions called the *exponential family* provides the right balance of generality while retaining favourable computation and inference properties. Before we introduce the exponential family, let us see three more members of “named” probability distributions, the Bernoulli (Example 6.8), Binomial (Example 6.9) and Beta (Example 6.10) distributions.

“Computers” were a job description.

exponential family

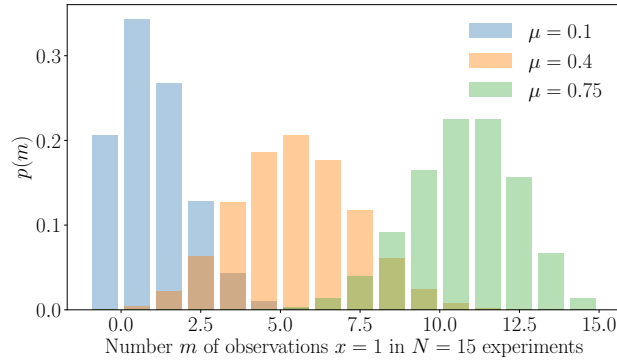
Example 6.8

The *Bernoulli distribution* is a distribution for a single binary variable $x \in \{0, 1\}$ and is governed by a single continuous parameter $\mu \in [0, 1]$ that represents the probability of $x = 1$. The Bernoulli distribution is defined as

Bernoulli distribution

$$p(x | \mu) = \mu^x (1 - \mu)^{1-x}, \quad x \in \{0, 1\}, \quad (6.95)$$

Figure 6.10
Examples of the
Binomial
distribution for
 $\mu \in \{0.1, 0.4, 0.75\}$
and $N = 15$.



$$\mathbb{E}[x] = \mu, \quad (6.96)$$

$$\mathbb{V}[x] = \mu(1 - \mu), \quad (6.97)$$

where $\mathbb{E}[x]$ and $\mathbb{V}[x]$ are the mean and variance of the binary random variable x .

An example where the Bernoulli distribution can be used is when we are interested in modeling the probability of “head” when flipping a coin.

Remark. The rewriting above of the Bernoulli distribution, where we use Boolean variables as numerical 0 or 1 and express them in the exponents, is a trick that is often used in machine learning textbooks. Another occurrence of this is when expressing the Multinomial distribution. \diamond

Binomial
distribution

Example 6.9

The *Binomial distribution* is a generalization of the Bernoulli distribution to a distribution over integers. In particular, the Binomial can be used to describe the probability of observing m occurrences of $x = 1$ in a set of N samples from a Bernoulli distribution where $p(x = 1) = \mu \in [0, 1]$. The Binomial distribution is defined as

$$p(m | N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}, \quad (6.98)$$

$$\mathbb{E}[m] = N\mu, \quad (6.99)$$

$$\mathbb{V}[m] = N\mu(1 - \mu) \quad (6.100)$$

where $\mathbb{E}[m]$ and $\mathbb{V}[m]$ are the mean and variance of m , respectively.

An example where the Binomial could be used is if we want to describe the probability of observing m “heads” in N coin-flip experiments if the probability for observing head in a single experiment is μ .

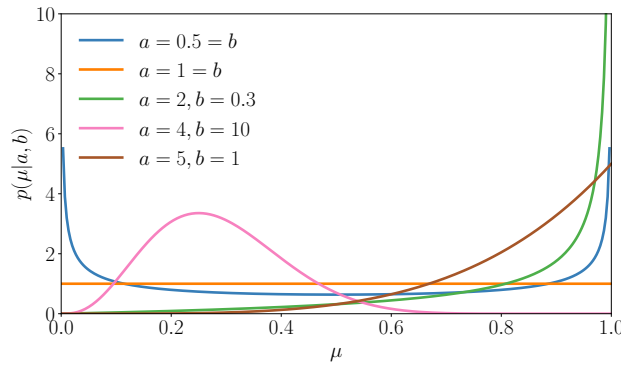


Figure 6.11
Examples of the
Beta distribution for
different values of α
and β .

Example 6.10

We may wish to model a continuous random variable on a finite interval. The *Beta distribution* is a distribution over a continuous random variable $\mu \in [0, 1]$, which is often used to represent the probability for some binary event (e.g., the parameter governing the Bernoulli distribution). The Beta distribution (illustrated in Figure 6.11) itself is governed by two parameters $\alpha > 0$, $\beta > 0$ and is defined as

$$p(\mu | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha-1} (1 - \mu)^{\beta-1} \quad (6.101)$$

$$\mathbb{E}[\mu] = \frac{\alpha}{\alpha + \beta}, \quad \mathbb{V}[\mu] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} \quad (6.102)$$

where $\Gamma(\cdot)$ is the Gamma function defined as

$$\Gamma(t) := \int_0^\infty x^{t-1} \exp(-x) dx, \quad t > 0. \quad (6.103)$$

$$\Gamma(t + 1) = t\Gamma(t). \quad (6.104)$$

Note that the fraction of Gamma functions in (6.101) normalizes the Beta distribution.

Beta distribution

Intuitively, α moves probability mass toward 1, whereas β moves probability mass toward 0. There are some special cases (Murphy, 2012):

- For $\alpha = 1 = \beta$ we obtain the uniform distribution $\mathcal{U}[0, 1]$.
- For $\alpha, \beta < 1$, we get a bimodal distribution with spikes at 0 and 1.
- For $\alpha, \beta > 1$, the distribution is unimodal.
- For $\alpha, \beta > 1$ and $\alpha = \beta$, the distribution is unimodal, symmetric and centered in the interval $[0, 1]$, i.e., the mode/mean is at $\frac{1}{2}$.

Remark. There is a whole zoo of distributions with names, and they are related in different ways to each other (Leemis and McQueston, 2008). It is worth keeping in mind that each named distribution is created for a

particular reason, but may have other applications. Knowing the reason behind the creation of a particular distribution often allows insight into how to best use it. We introduced the above three distributions to be able to illustrate the concepts of conjugacy (Section 6.6.1) and exponential families (Section 6.6.3). \diamond

6.6.1 Conjugacy

According to Bayes' theorem (6.22), the posterior is proportional to the product of the prior and the likelihood. The specification of the prior can be tricky for two reasons: First, the prior should encapsulate our knowledge about the problem before we see some data. This is often difficult to describe. Second, it is often not possible to compute the posterior distribution analytically. However, there are some priors that are computationally convenient: *conjugate priors*.

conjugate priors

conjugate

Definition 6.13 (Conjugate Prior). A prior is *conjugate* for the likelihood function if the posterior is of the same form/type as the prior.

Conjugacy is particularly convenient because we can algebraically calculate our posterior distribution by updating the parameters of the prior distribution.

Remark. When considering the geometry of probability distributions, conjugate priors retain the same distance structure as the likelihood (Agarwal and III, 2010). \diamond

To introduce a concrete example of conjugate priors, we describe below the Binomial distribution (defined on discrete random variables) and the Beta distribution (defined on continuous random variables).

Example 6.11 (Beta-Binomial Conjugacy)

Consider a Binomial random variable $x \sim \text{Bin}(N, \mu)$ where

$$p(x | N, \mu) = \binom{N}{x} \mu^x (1 - \mu)^{N-x} \quad x = 0, 1, \dots, N \quad (6.105)$$

is the probability of finding x times the outcome “head” in N coin flips, where μ is the probability of a “head”. We place a Beta prior on the parameter μ , that is $\mu \sim \text{Beta}(\alpha, \beta)$ where

$$p(\mu | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha-1} (1 - \mu)^{\beta-1} \quad (6.106)$$

If we now observe some outcome $x = h$, that is we see h heads in N coin flips, we compute the posterior distribution on μ as

$$p(\mu | x = h, N, \alpha, \beta) \propto p(x | N, \mu) p(\mu | \alpha, \beta) \quad (6.107a)$$

Likelihood	Conjugate prior	Posterior
Bernoulli	Beta	Beta
Binomial	Beta	Beta
Gaussian	Gaussian/inverse Gamma	Gaussian/inverse Gamma
Gaussian	Gaussian/inverse Wishart	Gaussian/inverse Wishart
Multinomial	Dirichlet	Dirichlet

Table 6.2 Examples of conjugate priors for common likelihood functions.

$$= \mu^h (1 - \mu)^{(N-h)} \mu^{\alpha-1} (1 - \mu)^{\beta-1} \quad (6.107b)$$

$$= \mu^{h+\alpha-1} (1 - \mu)^{(N-h)+\beta-1} \quad (6.107c)$$

$$\propto \text{Beta}(h + \alpha, N - h + \beta) \quad (6.107d)$$

i.e., the posterior distribution is a Beta distribution as the prior, i.e., the Beta prior is conjugate for the parameter μ in the Binomial likelihood function.

In the following example, we will derive a result that is similar to the Beta-Binomial conjugacy result. Here we will show that the Beta distribution is a conjugate prior for the Bernoulli distribution.

Example 6.12 (Beta-Bernoulli Conjugacy)

Let $x \in \{0, 1\}$ be distributed according to the Bernoulli distribution with parameter $\theta \in [0, 1]$, that is $p(x = 1 | \theta) = \theta$. This can also be expressed as $p(x | \theta) = \theta^x (1 - \theta)^{1-x}$. Let θ be distributed according to a Beta distribution with parameters α, β , that is $p(\theta | \alpha, \beta) \propto \theta^{\alpha-1} (1 - \theta)^{\beta-1}$.

Multiplying the Beta and the Bernoulli distributions, we get

$$p(\theta | x, \alpha, \beta) = p(x | \theta) \times p(\theta | \alpha, \beta) \quad (6.108a)$$

$$\propto \theta^x (1 - \theta)^{1-x} \times \theta^{\alpha-1} (1 - \theta)^{\beta-1} \quad (6.108b)$$

$$= \theta^{\alpha+x-1} (1 - \theta)^{\beta+(1-x)-1} \quad (6.108c)$$

$$\propto p(\theta | \alpha + x, \beta + (1 - x)). \quad (6.108d)$$

The last line above is the Beta distribution with parameters $(\alpha + x, \beta + (1 - x))$.

Table 6.2 lists examples for conjugate priors for the parameters of some standard likelihoods used in probabilistic modeling. Distributions such as Multinomial, inverse Gamma, inverse Wishart, and Dirichlet can be found in any statistical text, and are for example described in Bishop (2006).

The Beta distribution is the conjugate prior for the parameter μ in both the Binomial and the Bernoulli likelihood. For a Gaussian likelihood function, we can place a conjugate Gaussian prior on the mean. The reason why the Gaussian likelihood appears twice in the table is that we need to distinguish the univariate from the multivariate case. In the univariate

Alternatively, the
Gamma prior is
conjugate for the
precision (inverse
variance) in the
Gaussian likelihood.
Alternatively, the
Wishart prior is
conjugate for the
precision matrix
(inverse covariance
matrix) in the
Gaussian likelihood.

sufficient statistics

(scalar) case, the inverse Gamma is the conjugate prior for the variance. In the multivariate case, we use a conjugate inverse Wishart distribution as a prior on the covariance matrix. The Dirichlet distribution is the conjugate prior for the multinomial likelihood function. For further details, we refer to Bishop (2006).

6.6.2 Sufficient Statistics

Recall that a statistic of a random variable is a deterministic function of that random variable. For example if $\mathbf{x} = [x_1, \dots, x_N]^\top$ is a vector of univariate Gaussian random variables, that is $x_n \sim \mathcal{N}(\mu, \sigma^2)$, then the sample mean $\hat{\mu} = \frac{1}{N}(x_1 + \dots + x_N)$ is a statistic. Sir Ronald Fisher discovered the notion of *sufficient statistics*: the idea that there are statistics that will contain all available information that can be inferred from data corresponding to the distribution under consideration. In other words sufficient statistics carry all the information needed to make inference about the population, that is they are the statistics that are sufficient to represent the distribution.

For a set of distributions parameterized by θ , let x be a random variable with distribution given an unknown θ_0 . A vector $\phi(x)$ of statistics are called sufficient statistics for θ_0 if they contain all possible information about θ_0 . To be more formal about “contain all possible information”: this means that the probability of x given θ can be factored into a part that does not depend on θ , and a part that depends on θ only via $\phi(x)$. The Fisher-Neyman factorization theorem formalizes this notion, which we state below without proof.

Theorem 6.14 (Fisher-Neyman). *Let x have probability density function $p(x | \theta)$. Then the statistics $\phi(x)$ are sufficient for θ if and only if $p(x | \theta)$ can be written in the form*

$$p(x | \theta) = h(x)g_\theta(\phi(x)). \quad (6.109)$$

where $h(x)$ is a distribution independent of θ and g_θ captures all the dependence on θ via sufficient statistics $\phi(x)$.

If $p(x | \theta)$ does not depend on θ then $\phi(x)$ is trivially a sufficient statistic for any function ϕ . The more interesting case is that $p(x | \theta)$ is dependent only on $\phi(x)$ and not x itself. In this case, $\phi(x)$ is a sufficient statistic for x .

In machine learning we consider a finite number of samples from a distribution. One could imagine that for simple distributions (such as the Bernoulli in Example 6.8) we only need a small number of samples to estimate the parameters of the distributions. We could also consider the opposite problem: if we have a set of data (a sample from an unknown distribution), which distribution gives the best fit? A natural question to

ask is as we observe more data, do we need more parameters θ to describe the distribution? It turns out that the answer is yes in general, and this is studied in non-parametric statistics (Wasserman, 2007). A converse question is to consider which class of distributions have finite dimensional sufficient statistics, that is the number of parameters needed to describe them do not increase arbitrarily. The answer is exponential family distributions, described in the following section.

6.6.3 Exponential Family

There are three possible levels of abstraction we can have when considering distributions (of discrete or continuous random variables). At level one (the most concrete end of the spectrum), we have a particular named distribution with fixed parameters, for example a univariate Gaussian $\mathcal{N}(0, 1)$ with zero mean and unit variance. In machine learning we often use the second level of abstraction, that is we fix the parametric form (the univariate Gaussian) and infer the parameters from data. For example, we assume a univariate Gaussian $\mathcal{N}(\mu, \sigma^2)$ with unknown mean μ and unknown variance σ^2 , and use a maximum likelihood fit to determine the best parameters (μ, σ^2) . We will see an example of this when considering linear regression in Chapter 9. A third level of abstraction is to consider families of distributions, and in this book, we consider the exponential family. The univariate Gaussian is an example of a member of the exponential family. Many of the widely used statistical models, including all the “named” models in Table 6.2, are members of the exponential family. They can all be unified into one concept (Brown, 1986).

Remark. A brief historical anecdote: like many concepts in mathematics and science, exponential families were independently discovered at the same time by different researchers. In the years 1935–1936, Edwin Pitman in Tasmania, Georges Darmon in Paris, and Bernard Koopman in New York, independently showed that the exponential families are the only families that enjoy finite-dimensional sufficient statistics under repeated independent sampling (Lehmann and Casella, 1998). \diamond

An *exponential family* is a family of probability distributions, parameterized by $\theta \in \mathbb{R}^D$, of the form

exponential family

$$p(\mathbf{x} | \theta) = h(\mathbf{x}) \exp(\langle \theta, \phi(\mathbf{x}) \rangle - A(\theta)) , \quad (6.110)$$

where $\phi(\mathbf{x})$ is the vector of sufficient statistics. In general, any inner product (Section 3.2) can be used in (6.110), and for concreteness we will use the standard dot product here ($\langle \theta, \phi(\mathbf{x}) \rangle = \theta^\top \phi(\mathbf{x})$). Note that the form of the exponential family is essentially a particular expression of $g_\theta(\phi(\mathbf{x}))$ in the Fisher-Neyman theorem (Theorem 6.14).

The factor $h(\mathbf{x})$ can be absorbed into the dot product term by adding another entry to the vector of sufficient statistics $\log h(\mathbf{x})$, and constrain-

log partition
function

ing the corresponding parameter $\theta_0 = 1$. The term $A(\boldsymbol{\theta})$ is the normalization constant that ensures that the distribution sums up or integrates to one and is called the *log partition function*. A good intuitive notion of exponential families can be obtained by ignoring these two terms and considering exponential families as distributions of the form

$$p(\mathbf{x} | \boldsymbol{\theta}) \propto \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x})). \quad (6.111)$$

natural parameters

3826

3827

3828

3829

3830

3831

For this form of parameterization, the parameters $\boldsymbol{\theta}$ are called the *natural parameters*. At first glance it seems that exponential families is a mundane transformation by adding the exponential function to the result of a dot product. However, there are many implications that allow for convenient modelling and efficient computation based on the fact that we can capture information about data in $\boldsymbol{\phi}(\mathbf{x})$.

Example 6.13 (Gaussian as Exponential Family)

Consider the univariate Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$. Let $\boldsymbol{\phi}(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$. Then by using the definition of the exponential family,

$$p(x | \boldsymbol{\theta}) \propto \exp(\theta_1 x + \theta_2 x^2). \quad (6.112)$$

Setting

$$\boldsymbol{\theta} = \left[\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right]^\top \quad (6.113)$$

and substituting into (6.112) we obtain

$$p(x | \boldsymbol{\theta}) \propto \exp\left(\frac{\mu x}{\sigma^2} - \frac{x^2}{2\sigma^2}\right) \propto \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right). \quad (6.114)$$

Therefore, the univariate Gaussian distribution is a member of the exponential family with sufficient statistic $\boldsymbol{\phi}(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$.

Example 6.14 (Bernoulli as Exponential Family)

Recall the Bernoulli distribution from Example 6.8

$$p(x | \mu) = \mu^x (1 - \mu)^{1-x}, \quad x \in \{0, 1\}. \quad (6.115)$$

This can be written in exponential family form

$$p(x | \mu) = \exp[\log(\mu^x (1 - \mu)^{1-x})] \quad (6.116)$$

$$= \exp[x \log \mu + (1 - x) \log(1 - \mu)] \quad (6.117)$$

$$= \exp[x \log \mu - x \log(1 - \mu) + \log(1 - \mu)] \quad (6.118)$$

$$= \exp\left[x \log \frac{\mu}{1 - \mu} + \log(1 - \mu)\right]. \quad (6.119)$$

The last line (6.119) can be identified as being in exponential family form (6.110) by observing that

$$h(x) = 1 \quad (6.120)$$

$$\theta = \log \frac{\mu}{1 - \mu} \quad (6.121)$$

$$\phi(x) = x \quad (6.122)$$

$$A(\theta) = -\log(1 - \mu) = \log(1 + \exp(\theta)). \quad (6.123)$$

The relationship between θ and μ is invertible,

$$\mu = \frac{1}{1 + \exp(-\theta)}. \quad (6.124)$$

This relation is used to obtain the right equality of (6.123).

Remark. The relationship between the original Bernoulli parameter μ and the natural parameter θ is known as the *sigmoid* or logistic function. Observe that $\mu \in (0, 1)$ but $\theta \in \mathbb{R}$, and therefore the sigmoid function squeezes a real value into the range $(0, 1)$. This property is useful in machine learning, for example it is used in logistic regression (Bishop, 1995, Section 4.3.2), as well as as a nonlinear activation functions in neural networks (Goodfellow et al., 2016, Chapter 6). \diamond

sigmoid

It is often not obvious how to find the parametric form of the conjugate distribution of a particular distribution. Exponential families provide a convenient way to find conjugate pairs of distributions. Consider the random variable x distributed as an exponential family (6.110)

$$p(x | \theta) = h(x) \exp(\langle \theta, \phi(x) \rangle - A(\theta)). \quad (6.125)$$

Every exponential family has a conjugate prior (Brown, 1986)

$$p(\theta | \gamma) = h_c(\theta) \exp\left(\left\langle \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix}, \begin{bmatrix} \theta \\ -A(\theta) \end{bmatrix} \right\rangle - A_c(\gamma)\right), \quad (6.126)$$

where $\gamma = \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix}$ has dimension $\dim(\theta) + 1$. The sufficient statistics of the conjugate prior are $\begin{bmatrix} \theta \\ -A(\theta) \end{bmatrix}$. By using the knowledge of the general form of conjugate priors for exponential families, we can derive functional forms of conjugate priors corresponding to particular distributions.

Example 6.15

Recall the exponential family form of the Bernoulli distribution (6.119),

$$p(x | \mu) = \exp\left[x \log \frac{\mu}{1 - \mu} + \log(1 - \mu)\right]. \quad (6.127)$$

The canonical conjugate prior therefore has the same form

$$p(\mu | \gamma, n_0) = \exp \left[n_0 \gamma \log \frac{\mu}{1 - \mu} + n_0 \log(1 - \mu) - A_c(\gamma, n_0) \right], \quad (6.128)$$

which simplifies to

$$p(\mu | \gamma, n_0) = \exp [n_0 \gamma \log \mu + n_0(1 - \gamma) \log(1 - \mu) - A_c(\gamma, n_0)]. \quad (6.129)$$

Putting this in non-exponential family form

$$p(\mu | \gamma, n_0) \propto \mu^{n_0 \gamma} (1 - \mu)^{n_0(1 - \gamma)} \quad (6.130)$$

which is of the same form as the Beta distribution (6.101), with minor manipulations to get the original parametrization (Example 6.12).

Observe that in this example we have derived the form of the Beta distribution by looking at the conjugate prior of the exponential family.

As mentioned in the previous section, the main motivation for exponential families is that they have finite-dimensional sufficient statistics. Additionally, conjugate distributions are easy to write down, and the conjugate distributions also come from an exponential family. From an inference perspective, maximum likelihood estimation behaves nicely because empirical estimates of sufficient statistics are optimal estimates of the population values of sufficient statistics (recall the mean and covariance of a Gaussian). From an optimization perspective, the log-likelihood function is concave allowing for efficient optimization approaches to be applied (Chapter 7).

6.7 Change of Variables/Inverse Transform

It may seem that there are very many known distributions, but in reality the set of distributions for which we have names is quite limited. Therefore, it is often useful to understand how transformed random variables are distributed. For example, assume that x is a random variable distributed according to the univariate normal distribution $\mathcal{N}(0, 1)$, what is the distribution of x^2 ? Another example, which is quite common in machine learning, is: given that x_1 and x_2 are univariate standard normal, what is the distribution of $\frac{1}{2}(x_1 + x_2)$?

One option to work out the distribution of $\frac{1}{2}(x_1 + x_2)$ is to calculate the mean and variance of x_1 and x_2 and then combine them. As we saw in Section 6.4.4, we can calculate the mean and variance of resulting random variables when we consider affine transformations of random variables. However, we may not be able to obtain the functional form of the distribution under transformations. Furthermore, we may be interested

in nonlinear transformations of random variables for which closed-form expressions are not readily available.

Remark (Notation). In this section, we will be explicit about random variables and the values they take. Hence, we will use small letters x, y to denote random variables and small capital letters x, y to denote the values that the random variables take. We will explicitly write probability mass functions (pmf) of discrete random variables x as $P(x = x)$. For continuous random variables x , the probability density function (pdf) is written as $f(x)$ and the cumulative distribution function (cdf) is written as $F_x(x \leq x)$. \diamond

We will look at two approaches for obtaining distributions of transformations of random variables: a direct approach using the definition of a cumulative distribution function and a change-of-variable approach that uses the chain rule of calculus (Section 5.2.2). The change-of-variable approach is widely used because it provides a “recipe” for attempting to compute the resulting distribution due to a transformation. We will explain the techniques for univariate random variables, and will only briefly provide the results for the general case of multivariate random variables.

As mentioned in Section 6.1, random variables and probability distributions are closely associated with each other. It is worth carefully teasing apart the two ideas, and in doing so we will motivate why we need to transform random variables.

Moment generating functions can also be used to study transformations of random variables (Casella and Berger, 2002, Chapter 2).

Example 6.16

Consider a medical test that returns the number of cancerous cells that can be found in the biopsy. The state space is the set of non-negative integers. The random variable x is the *square* of the number of cancerous cells. Given that we know the probability distribution corresponding to the number of cancerous cells in a biopsy, how do we obtain the distribution of random variable x ?

Transformations of discrete random variables can be understood directly. Given a discrete random variable x with probability mass function (pmf) $P(x = x)$ (Section 6.2.1), and an invertible function $U(x)$. Consider the transformed random variable $y := U(x)$, with pmf $P(y = y)$. Then

$$P(y = y) = P(U(x) = y) \quad \text{transformation of interest} \quad (6.131a)$$

$$= P(x = U^{-1}(y)) \quad \text{inverse} \quad (6.131b)$$

where we can observe that $x = U^{-1}(y)$. Therefore for discrete random variables, transformations directly change the individual events (with the probabilities appropriately transformed).

6.7.1 Distribution Function Technique

3893

3894 The distribution function technique goes back to first principles, and uses
 3895 the definition of a cumulative distribution function (cdf) $F_x(x) = P(x \leq$
 3896 $x)$ and the fact that its differential is the probability density function (pdf)
 3897 $f(x)$ (Wasserman, 2004, Chapter 2). For a random variable x and a func-
 3898 tion U , we find the pdf of the random variable $y := U(x)$ by

1 Finding the cdf:

$$F_y(Y) = P(y \leq Y) \quad (6.132)$$

2 Differentiating the cdf $F_y(Y)$ to get the pdf $f(y)$.

$$f(y) = \frac{d}{dy} F_y(Y). \quad (6.133)$$

3899 We also need to keep in mind that the domain of the random variable may
 3900 have changed due to the transformation by U .

Example 6.17

Let x be a continuous random variable with probability density function on $0 \leq x \leq 1$

$$f(x) = 3x^2. \quad (6.134)$$

We are interested in finding the pdf of $y = x^2$.

The function f is an increasing function of x , and the resulting value of y lies in the interval $[0, 1]$. We obtain

$$F_y(Y) = P(y \leq Y) \quad \text{definition of cdf} \quad (6.135a)$$

$$= P(x^2 \leq Y) \quad \text{transformation of interest} \quad (6.135b)$$

$$= P(x \leq Y^{\frac{1}{2}}) \quad \text{inverse} \quad (6.135c)$$

$$= F_x(Y^{\frac{1}{2}}) \quad \text{definition of cdf} \quad (6.135d)$$

$$= \int_0^{Y^{\frac{1}{2}}} 3t^2 dt \quad \text{cdf as a definite integral} \quad (6.135e)$$

$$= [t^3]_{t=0}^{t=Y^{\frac{1}{2}}} \quad \text{result of integration} \quad (6.135f)$$

$$= Y^{\frac{3}{2}}, \quad 0 \leq Y \leq 1. \quad (6.135g)$$

Therefore, the cdf of y is

$$F_y(Y) = Y^{\frac{3}{2}} \quad (6.136)$$

for $0 \leq Y \leq 1$. To obtain the pdf, we differentiate the cdf

$$f(y) = \frac{d}{dy} F_y(Y) = \frac{3}{2} y^{\frac{1}{2}} \quad (6.137)$$

for $0 \leq y \leq 1$.

In Example 6.17, we considered a strictly monotonically increasing function $f(x) = 3x^2$. This means that we could compute an inverse function. In general, we require that the function of interest $y = U(x)$ has an inverse $x = U^{-1}(y)$. A useful result can be obtained by considering the cumulative distribution function $F_x(x)$ of a random variable x , and using it as the transformation $U(x)$. This leads to the following theorem which is called the probability integral transform. The result is the basis of generating samples from distributions whose cdfs are known by first generating a sample from a uniform distribution and then transforming it by the inverse cdf.

Functions that have inverses are called injective functions (Section 2.7).

Theorem 6.15. *This is Theorem 2.1.10 in Casella and Berger (2002). Let x be a continuous random variable with a strictly monotonic cumulative distribution function $F_x(\cdot)$. Then the random variable y defined as*

$$y = F_x(x), \quad (6.138)$$

has a uniform distribution.

Proof We need to show that the cumulative distribution function of y defines a distribution of a uniform random variable. Recall that by the axioms of probability (Section 6.1) probabilities must be non-negative and sum/integrate to one. Therefore, the range of possible values of $y = F_x(x)$ is the interval $[0, 1]$. For any $F_x(\cdot)$, the inverse $F_x^{-1}(\cdot)$ exists because we assumed that $F_x(\cdot)$ is strictly monotonically increasing, which we will use in the following.

Given any continuous random variable x , the definition of a cdf gives

$$F_y(y) = P(y \leq Y) \quad (6.139a)$$

$$= P(F_x(x) \leq Y) \quad \text{transformation of interest} \quad (6.139b)$$

$$= P(x \leq F_x^{-1}(Y)) \quad \text{inverse exists} \quad (6.139c)$$

$$= F_x(F_x^{-1}(Y)) \quad \text{definition of cdf} \quad (6.139d)$$

$$= Y, \quad (6.139e)$$

where the last line is due to the fact that $F_x(\cdot)$ composed with its inverse results in an identity transformation. The statement $F_y(y) = y$ along with the fact that y lies in the interval $[0, 1]$ means that $F_y(\cdot)$ is the cdf of the uniform random variable on the unit interval. \square

Theorem 6.15 is known as the *probability integral transform*, and it is used to derive algorithms for sampling from distributions by transforming the result of sampling from a uniform random variable (Bishop, 2006). It is also used for hypothesis testing whether a sample comes from a particular distribution (Lehmann and Romano, 2005). The idea that the output of a cdf gives a uniform distribution also forms the basis of copulas (Nelsen, 2006).

probability integral transform

6.7.2 Change of Variables

The distribution function technique in Section 6.7.1 is derived from first principles, based on the definitions of cdfs and using properties of inverses, differentiation and integration. This argument from first principles relies on two facts:

- 1 We can transform the cdf of y into an expression that is a cdf of x .
- 2 We can differentiate the cdf to obtain the pdf.

Let us break down the reasoning step by step, with the goal of understanding the more general change of variables approach in Theorem 6.16.

Remark. The name change of variables comes from the idea of changing the variable of integration when faced with a difficult integral. For univariate functions, we use the substitution rule of integration,

$$\int f(g(x))g'(x)dx = \int f(u)du \quad \text{where } u = g(x). \quad (6.140)$$

The derivation of this rule is based on the chain rule of calculus (5.31) and by applying twice the fundamental theorem of calculus. The fundamental theorem of calculus formalizes the fact that integration and differentiation are somehow “inverses” of each other. An intuitive understanding of the rule can be obtained by thinking (loosely) about small changes (differentials) to the equation $u = g(x)$. That is by considering $\Delta u = g'(x)\Delta x$ as a differential of $u = g(x)$. By substituting $u = g(x)$, the argument inside the integral on the right hand side of (6.140) becomes $f(g(x))$. By pretending that the term du can be approximated by $du \approx \Delta u = g'(x)\Delta x$, and that $dx \approx \Delta x$, we obtain (6.140). \diamond

Consider a function of a random variable $y = U(x)$, where $x \in [a, b]$. By the definition of the cdf, we have

$$F_y(y) = P(y \leq Y). \quad (6.141)$$

We are interested in a function U of the random variable

$$P(y \leq Y) = P(U(x) \leq Y), \quad (6.142)$$

where we assume that the function U is invertible. By applying the inverse U^{-1} to the arguments of $P(U(x) \leq Y)$, we obtain

$$P(U(x) \leq Y) = P(U^{-1}(U(x)) \leq U^{-1}(Y)) = P(x \leq U^{-1}(Y)), \quad (6.143)$$

which is an expression of the cdf of x . Recall the definition of the cdf in terms of the pdf

$$P(x \leq U^{-1}(Y)) = \int_a^{U^{-1}(Y)} f(x)dx. \quad (6.144)$$

Change of variables in probability relies on the change of variables method in calculus (Tandra, 2014).

Now we have an expression of the cdf of y in terms of x :

$$F_y(Y) = \int_a^{U^{-1}(Y)} f(x) dx. \quad (6.145)$$

To obtain the pdf, we differentiate (6.145) with respect to y .

$$f(y) = \frac{d}{dy} F_y(Y) = \frac{d}{dy} \int_a^{U^{-1}(Y)} f(x) dx \quad (6.146)$$

Note that the integral on the right hand side is with respect to x , but we need an integral with respect to y because we are differentiating with respect to y . In particular we use (6.140) to get the substitution

$$\int f(U^{-1}(y)) U^{-1'}(y) dy = \int f(x) dx \quad \text{where } x = U^{-1}(y). \quad (6.147)$$

Using (6.147) on the right hand side of (6.146) gives us

$$f(y) = \frac{d}{dy} \int_a^{U^{-1}(Y)} f_x(U^{-1}(y)) U^{-1'}(y) dy. \quad (6.148)$$

We then recall that differentiation is a linear operator and we use the subscript x to remind ourselves that $f_x(U^{-1}(y))$ is a function of x and not y . Invoking the fundamental theorem of calculus again gives us

$$f(y) = f_x(U^{-1}(y)) \times \left| \frac{d}{dy} U^{-1}(y) \right|. \quad (6.149)$$

This is called the *change of variable* technique. The term $\left| \frac{d}{dy} U^{-1}(y) \right|$ measures how much a unit volume changes when applying U (see also the Remark on page 151. In (6.149) we introduced the absolute value of the differential. For decreasing functions, it turns out that an additional negative sign is needed, and instead of having two types of change-of-variable rules, the absolute value unifies both of them.

Remark. In comparison to the discrete case in (6.131b), we have an additional factor $\left| \frac{d}{dy} U^{-1}(y) \right|$. The continuous case requires more care because $P(y = Y) = 0$ for all Y . The probability density function $f(y)$ does not have a description as a probability of an event involving y . \diamond

So far in this section we have been studying univariate change of variables. The case for multivariate random variables is analogous, but complicated by fact that the absolute value cannot be used for multivariate functions. Instead we use the determinant of the Jacobian matrix. Recall from (5.56) that the Jacobian is a matrix of partial derivatives, and that the existence of a non-zero determinant shows that we can invert the Jacobian. Recall the discussion in Section 4.1 that the determinant arises because our differentials (cubes of volume) are transformed into parallelepipeds by the Jacobian. Let us summarize the discussion above in the

3968 following theorem, which gives us a recipe for multivariate change of vari-
 3969 ables.

Theorem 6.16. [Theorem 17.2 in Billingsley (1995)] Let $f(\mathbf{x})$ be the value of the probability density of the multivariate continuous random variable \mathbf{x} . If the vector-valued function $\mathbf{y} = U(\mathbf{x})$ is differentiable and invertible for all values within the domain of \mathbf{x} , then for corresponding values of \mathbf{y} , the probability density of $\mathbf{y} = U(\mathbf{x})$ is given by

$$f(\mathbf{y}) = f_{\mathbf{x}}(U^{-1}(\mathbf{y})) \times \left| \det \left(\frac{\partial}{\partial \mathbf{y}} U^{-1}(\mathbf{y}) \right) \right|. \quad (6.150)$$

3970 The theorem looks intimidating at first glance, but the key point is that
 3971 a change of variable of a multivariate random variable follows the pro-
 3972 cedure of the univariate change of variable. First we need to work out
 3973 the inverse transform, and substitute that into the density of \mathbf{x} . Then we
 3974 calculate the determinant of the Jacobian and multiply the result. The
 3975 following example illustrates the case of a bivariate random variable.

Example 6.18

Consider a bivariate random variable $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ with probability density function

$$f \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \frac{1}{2\pi} \exp \left(-\frac{1}{2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right). \quad (6.151)$$

We use the change-of-variable technique from Theorem 6.16 to derive the effect of a linear transformation (Section 2.7) of the random variable. Consider a matrix $\mathbf{A} \in \mathbb{R}^{2 \times 2}$ defined as

$$\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}. \quad (6.152)$$

We are interested in finding the probability density function of the transformed bivariate random variable $\mathbf{y} = \mathbf{A}\mathbf{x}$.

Recall that for change of variables we require the inverse transformation of \mathbf{x} as a function of \mathbf{y} . Since we consider linear transformations, the inverse transformation is given by the matrix inverse (see Section 2.2.2). For 2×2 matrices, we can explicitly write out the formula, given by

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{A}^{-1} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}. \quad (6.153)$$

Observe that $ad - bc$ is the determinant (Section 4.1) of \mathbf{A} . The corresponding probability density function is given by

$$f(\mathbf{x}) = f(\mathbf{A}^{-1}\mathbf{y}) = \frac{1}{2\pi} \exp \left(-\frac{1}{2} \mathbf{y}^\top \mathbf{A}^{-\top} \mathbf{A}^{-1} \mathbf{y} \right). \quad (6.154)$$

The partial derivative of a matrix times a vector with respect to the vector is the matrix itself (Section 5.5) and, therefore,

$$\frac{\partial}{\partial \mathbf{y}} \mathbf{A}^{-1} \mathbf{y} = \mathbf{A}^{-1}. \quad (6.155)$$

Recall from Section 4.1 that the determinant of the inverse is the inverse of the determinant so that the determinant of the Jacobian matrix is given by

$$\det \left(\frac{\partial}{\partial \mathbf{y}} \mathbf{A}^{-1} \mathbf{y} \right) = \frac{1}{ad - bc}. \quad (6.156)$$

We are now able to apply the change-of-variable formula from Theorem 6.16 by multiplying (6.154) with (6.156), which yields

$$f(\mathbf{y}) = f(\mathbf{x}) \left| \det \left(\frac{\partial}{\partial \mathbf{y}} \mathbf{A}^{-1} \mathbf{y} \right) \right| \quad (6.157a)$$

$$= \frac{1}{2\pi} \exp \left(-\frac{1}{2} \mathbf{y}^\top \mathbf{A}^{-\top} \mathbf{A}^{-1} \mathbf{y} \right) |ad - bc|^{-1}. \quad (6.157b)$$

While Example 6.18 is based on a bivariate random variable, which allows to easily compute the matrix inverse, the relation above holds for higher dimensions.

Remark. We saw in Section 6.5 that the density $f(\mathbf{x})$ above is actually the standard Gaussian distribution, and the transformed density $f(\mathbf{y})$ is a bivariate Gaussian with covariance $\Sigma = \mathbf{A} \mathbf{A}^\top$. \diamond

6.8 Further Reading

This chapter is rather terse at times, Grinstead and Snell (1997) and Walpole et al. (2011) provides more relaxed presentations that are suitable for self study. Readers interested in more philosophical aspects of probability should consider Hacking (2001), whereas a more software engineering approach is presented by Downey (2014). An overview of exponential families can be found in Barndorff-Nielsen (2014). We will see more about how to use probability distributions to model machine learning tasks in Chapter 8. Ironically the recent surge in interest in neural networks has resulted in a broader appreciation of probabilistic models. For example the idea of normalizing flows (Rezende and Mohamed, 2015) relies on change of variables for transforming random variables. An overview of methods for variational inference as applied to neural networks is described in Chapters 16 to 20 of Goodfellow et al. (2016).

We side stepped a large part of the difficulty in continuous random variables by avoiding measure theoretic questions (Billingsley, 1995; Pollard, 2002), and by assuming without construction that we have real numbers,

and ways of defining sets on real numbers as well as their appropriate frequency of occurrence. These details do matter, for example in the specification of conditional probability $p(y|x)$ for continuous random variables x, y (Proschan and Presnell, 1998). The lazy notation hides the fact that we want to specify that $x = \mathbf{x}$ (which is a set of measure zero). Furthermore we are interested in the probability density function of y . A more precise notation would have to say $\mathbb{E}_y[f(y) | \sigma(x)]$, where we take the expectation over y of a test function f conditioned on the σ -algebra of x . A more technical audience interested in the details of probability theory have many options (Jacod and Protter, 2004; Jaynes, 2003; MacKay, 2003b; Grimmett and Welsh, 2014) including some very technical discussions (Çinlar, 2011; Dudley, 2002; Shirayev, 1984; Lehmann and Casella, 1998; Bickel and Doksum, 2006). As machine learning allows us to model more intricate distributions on ever more complex types of data, a developer of probabilistic machine learning models would have to understand these more technical aspects. Machine learning books with a probabilistic modelling focus includes MacKay (2003b); Bishop (2006); Murphy (2012); Barber (2012); Rasmussen and Williams (2006).

Exercises

6.1 Consider a mixture of two Gaussian distributions (illustrated in Figure 6.4)

$$\mathcal{N}\left(\begin{bmatrix} 10 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) + 1.5\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 & 2.9 \\ -1.1 & 0.7 \end{bmatrix}\right).$$

- 1 Compute the marginal distributions for each dimension
 - 2 Compute the mean, mode and median for each marginal distribution
 - 3 Compute the mean and mode for the 2 dimensional distribution
- 6.2 You have written a computer program that sometimes compiles and sometimes not (code does not change). You decide to model the apparent stochasticity (success vs no success) x of the compiler using a Bernoulli distribution with parameter μ :

$$p(x|\mu) = \mu^x(1-\mu)^{1-x}, \quad x \in \{0, 1\}$$

Choose a conjugate prior for the Bernoulli likelihood and compute the posterior distribution $p(\mu|x_1, \dots, x_N)$.

6.3 Consider the following time-series model:

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{A}\mathbf{x}_t + \mathbf{w}, & \mathbf{w} &\sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \\ \mathbf{y}_t &= \mathbf{C}\mathbf{x}_t + \mathbf{v}, & \mathbf{v} &\sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \end{aligned}$$

where \mathbf{w}, \mathbf{v} are i.i.d. Gaussian noise variables. Further, assume that $p(\mathbf{x}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$.

- 1 What is the form of $p(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T)$? Justify your answer (you do not have to explicitly compute the joint distribution). (1–2 sentences)
- 2 Assume that $p(\mathbf{x}_t|\mathbf{y}_1, \dots, \mathbf{y}_t) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$.

- 4028 a) Compute $p(\mathbf{x}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_t)$
- 4029 b) Compute $p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_t)$
- 4030 c) At time $t+1$, we observe the value $\mathbf{y}_{t+1} = \hat{\mathbf{y}}$. Compute $p(\mathbf{x}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_{t+1})$.
- 4031 6.4 Prove the relationship in Equation 6.44, which relates the standard defini-
4032 tion of the variance to the raw score expression for the variance.
- 4033 6.5 Prove the relationship in Equation 6.45, which relates the pairwise differ-
4034 ence between examples in a dataset with the raw score expression for the
4035 variance.
- 4036 6.6 Express the Bernoulli distribution in the natural parameter form of the ex-
4037ponential family (Equation (6.110)).
- 4038 6.7 Express the Binomial distribution as an exponential family distribution. Also
4039 express the Beta distribution is an exponential family distribution. Show that
4040 the product of the Beta and the Binomial distribution is also a member of
4041 the exponential family.
- 4042 6.8 Derive the relationship in Section 6.5.2 in two ways:
- 4043 1 By completing the square
- 4044 2 By expressing the Gaussian in its exponential family form
- The *product* of two Gaussians $\mathcal{N}(\mathbf{x}|\mathbf{a}, \mathbf{A})\mathcal{N}(\mathbf{x}|\mathbf{b}, \mathbf{B})$ is an unnormalized Gaussian distribution $c\mathcal{N}(\mathbf{x}|\mathbf{c}, \mathbf{C})$ with

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \quad (6.158)$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b}) \quad (6.159)$$

$$c = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1}(\mathbf{a} - \mathbf{b})\right). \quad (6.160)$$

- 4045 Note that the normalizing constant c itself can be considered a (normalized)
4046 Gaussian distribution either in \mathbf{a} or in \mathbf{b} with an “inflated” covariance matrix
4047 $\mathbf{A} + \mathbf{B}$, i.e., $c = \mathcal{N}(\mathbf{a}|\mathbf{b}, \mathbf{A} + \mathbf{B}) = \mathcal{N}(\mathbf{b}|\mathbf{a}, \mathbf{A} + \mathbf{B})$.
- 6.9 **Iterated Expectations.**
Consider two random variables x, y with joint distribution $p(x, y)$. Show that:

$$\mathbb{E}_x[x] = \mathbb{E}_y[\mathbb{E}_x[x|y]]$$

- 4048 Here, $\mathbb{E}_x[x|y]$ denotes the expected value of x under the conditional distri-
4049bution $p(x|y)$.
- 6.10 **Manipulation of Gaussian Random Variables.**
Consider a Gaussian random variable $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$, where $\mathbf{x} \in \mathbb{R}^D$. Furthermore, we have

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{w}, \quad (6.161)$$

- 4050 where $\mathbf{y} \in \mathbb{R}^E$, $\mathbf{A} \in \mathbb{R}^{E \times D}$, $\mathbf{b} \in \mathbb{R}^E$, and $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{Q})$ is indepen-
4051 dent Gaussian noise. “Independent” implies that \mathbf{x} and \mathbf{w} are independent
4052 random variables and that \mathbf{Q} is diagonal.
- 4053 1 Write down the likelihood $p(\mathbf{y}|\mathbf{x})$.

4054 2 The distribution $p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x}$ is Gaussian.¹ Compute the mean
4055 $\boldsymbol{\mu}_y$ and the covariance $\boldsymbol{\Sigma}_y$. Derive your result in detail.

3 The random variable \mathbf{y} is being transformed according to the measurement mapping

$$\mathbf{z} = \mathbf{C}\mathbf{y} + \mathbf{v}, \quad (6.162)$$

4056 where $\mathbf{z} \in \mathbb{R}^F$, $\mathbf{C} \in \mathbb{R}^{F \times E}$, and $\mathbf{v} \sim \mathcal{N}(\mathbf{v} | \mathbf{0}, \mathbf{R})$ is independent Gaussian
4057 (measurement) noise.

- 4058 • Write down $p(\mathbf{z}|\mathbf{y})$.
- 4059 • Compute $p(\mathbf{z})$, i.e., the mean $\boldsymbol{\mu}_z$ and the covariance $\boldsymbol{\Sigma}_z$. Derive your
4060 result in detail.

4061 4 Now, a value $\hat{\mathbf{y}}$ is measured. Compute the posterior distribution $p(\mathbf{x}|\hat{\mathbf{y}})$.²
4062 *Hint for solution:* Start by explicitly computing the joint Gaussian $p(\mathbf{x}, \mathbf{y})$.

4063 This also requires to compute the cross-covariances $\text{Cov}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}, \mathbf{y}]$ and $\text{Cov}_{\mathbf{y}, \mathbf{x}}[\mathbf{y}, \mathbf{x}]$.
4064 Then, apply the rules for Gaussian conditioning.

4065 6.11 Probability integral transformation

4066 Given a continuous random variable x , with cdf $F_x(x)$. Show that the random
4067 variable $y = F_x(x)$ is uniformly distributed.

¹An affine transformation of the Gaussian random variable \mathbf{x} into $\mathbf{A}\mathbf{x} + \mathbf{b}$ preserves Gaussianity. Furthermore, the sum of this Gaussian random variable and the independent Gaussian random variable \mathbf{w} is Gaussian.

²This posterior is also Gaussian, i.e., we need to determine only its mean and covariance matrix.