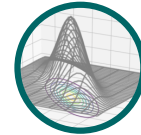


Probability and Distributions



3014 Probability, loosely speaking, concerns the study of uncertainty. Probabil-
 3015 ity can be thought of as the fraction of times an event occurs, or as a degree
 3016 of belief about an event. We then would like to use this probability to mea-
 3017 sure the chance of something occurring in an experiment. As mentioned
 3018 in Chapter 1, we often quantify: uncertainty in the data, uncertainty in the
 3019 machine learning model, and uncertainty in the predictions produced by
 3020 the model. Quantifying uncertainty requires the idea of a *random variable*,
 3021 which is a function that maps outcomes of random experiments to a set of
 3022 properties that we are interested in. Associated with the random variable
 3023 is a function which measures the probability that a particular outcome (or
 3024 set of outcomes) will occur, this is called the *probability distribution*.

random variable

probability
distribution

3025 Probability distributions are used as a building block for other con-
 3026 cepts, such as probabilistic modeling (Section 8.3), graphical models (Sec-
 3027 tion 8.4) and model selection (Section 8.5). In the next section, we present
 3028 the three concepts that define a probability space (the state space, the
 3029 events and the probability of an event) and how they are related to a
 3030 fourth concept called the random variable. The presentation is deliber-
 3031 ately slightly hand wavy since a rigorous presentation would occlude the
 3032 main idea. An outline of the concepts presented in this chapter are shown
 3033 in Figure 6.1.

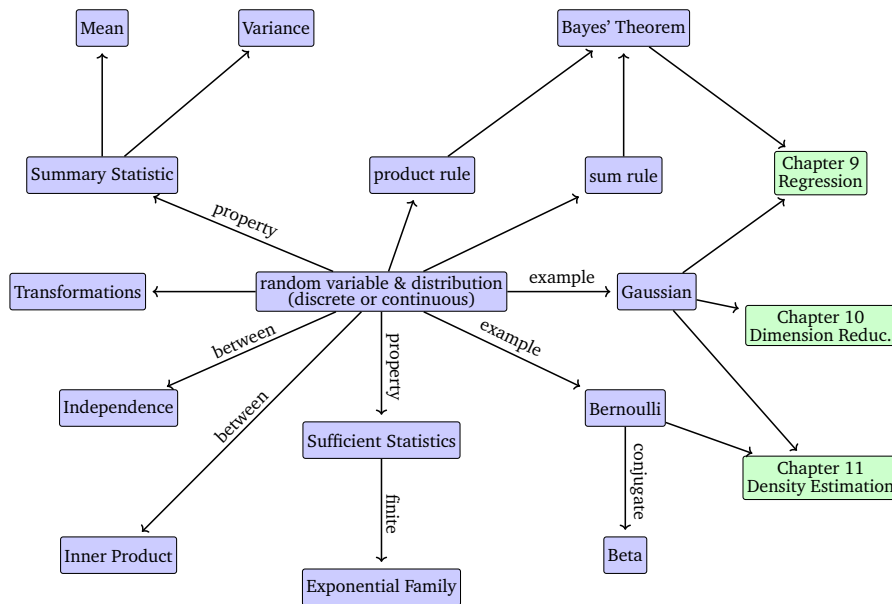
6.1 Construction of a Probability Space

3034
 3035 The theory of probability aims at defining a mathematical structure to
 3036 describe random outcomes of experiments. For example, when tossing a
 3037 single coin, one cannot determine the outcome, but by doing a large num-
 3038 ber of coin tosses, one can observe a regularity in the average outcome.
 3039 Using this mathematical structure of probability, the goal is to perform
 3040 automated reasoning, and in this sense probability generalizes logical rea-
 3041 soning (Jaynes, 2003).

6.1.1 Philosophical Issues

3042
 3043 When constructing automated reasoning systems, classical Boolean logic
 3044 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.



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).

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

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

3069 through different means, the same plausibility value must be found
3070 in all cases.

- 3071 b) Honesty: All available data must be taken into account.
3072 c) Reproducibility: If our state of knowledge about two problems are the
3073 same, then we must assign the same degree of plausibility to both of
3074 them.

3075 The Cox-Jaynes's theorem proves these plausibilities to be sufficient to
3076 define the universal mathematical rules that apply to plausibility p , up to
3077 transformation by an arbitrary monotonic function. Crucially, these rules
3078 are the rules of probability.

3079 *Remark.* In machine learning and statistics, there are two major interpre-
3080 tations of probability: the Bayesian and frequentist interpretations (Bishop,
3081 2006). The Bayesian interpretation uses probability to specify the degree
3082 of uncertainty that the user has about an event, and is sometimes referred
3083 to as subjective probability or degree of belief. The frequentist interpreta-
3084 tion considers probability to be the relative frequencies of events, in the
3085 limit when one has infinite data. \diamond

3086 Some machine learning texts on probabilistic models use lazy notation
3087 and jargon, which is confusing. Multiple distinct concepts are all referred
3088 to as “probability distribution”, and the reader has to often disentangle
3089 the meaning from the context. One trick to help make sense of probability
3090 distributions is to check whether we are trying to model something cate-
3091 gorical (a discrete random variable) or something continuous (a continuous
3092 random variable). The kinds of questions we tackle in machine learning
3093 are closely related to whether we are considering categorical or continu-
3094 ous models.

3095 6.1.2 Probability and Random Variables

3096 Modern probability is based on a set of axioms proposed by Kolmogorov (Ja-
3097 cod and Protter, 2004, Chapter 1 and 2) that introduce the three concepts
3098 of state space, event space and probability measure. The probability space
3099 models a real world process (referred to as an experiment) with random
3100 outcomes.

3101 The state space Ω

3102 The *state space* is the set of all possible outcomes of the experiment, state space
3103 usually denoted by Ω . For example, two successive coin tosses have a
3104 state space of {hh, tt, ht, th}, where “h” denotes “heads” and “t” denotes
3105 “tails”.

3106 The event space \mathcal{A}

3107 The *event space* is the space of potential results of the experiment. A event space
3108 subset A of the state space Ω is in the event space \mathcal{A} if at the end of the
3109 experiment we can observe whether a particular state $\omega \in \Omega$ is in A .

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 often
 3112 the 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 .

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
 3120 it to model some real world phenomenon. In machine learning we often
 3121 avoid explicitly referring to the probability space, but instead refer to
 3122 probabilities on quantities of interest which we denote by \mathcal{T} . We introduce
 3123 a function $x : \Omega \rightarrow \mathcal{T}$ that takes an element of Ω (an event) and returns
 3124 a particular quantity of interest, a value in \mathcal{T} . For example in the case of
 3125 tossing two coins and counting the number of heads, a random variable x
 3126 maps to the three possible outcomes: $x(\text{hh}) = 2$, $x(\text{ht}) = 1$, $x(\text{th}) = 1$ and
 3127 $x(\text{tt}) = 0$. In this particular case $\mathcal{T} = \{0, 1, 2\}$, and it is the probabilities
 3128 on elements of \mathcal{T} that we are interested in. This association or mapping
 3129 from Ω to \mathcal{T} is called a *random variable*. The name “random variable” is a
 3130 great source of misunderstanding as it is neither random nor is it a vari-
 3131 able. It is a function. For a finite state space Ω and finite \mathcal{T} , the function
 3132 corresponding to a random variable is essentially a look up table. For any
 3133 subset $S \subseteq \mathcal{T}$ we associate $P_x(S) \in [0, 1]$ (the probability) to a particu-
 3134 lar event occurring corresponding to the random variable x . Example 6.1
 3135 provides a concrete example illustrating the above terminology.

random variable

3136 *Remark.* The state space Ω above unfortunately is referred to by differ-
 3137 ent names in different books. Another common name for Ω is sample
 3138 space (Grinstead and Snell, 1997; Jaynes, 2003), and state space is some-
 3139 times reserved for referring to states in a dynamical system (Hasselblatt
 3140 and Katok, 2003). Other names sometimes used to describe Ω are: sample
 3141 description space, possibility space and (very confusingly) event space.
 3142 \diamond

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 $(\$, \$)$, $(\$, \mathcal{L})$, $(\mathcal{L}, \$)$, $(\mathcal{L}, \mathcal{L})$. 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 $S \subseteq \mathcal{T}$ (for example a single element of \mathcal{T} such as the outcome that one head is obtained when tossing two coins). 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 $S \subseteq \mathcal{T}$, we have the

following notation

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

3143 The left hand side of (6.8) is the probability of the set of possible outcomes
 3144 (e.g. number of heads = 1) that we are interested in. Via the random
 3145 variable x that maps states to outcomes, we see in the right hand side
 3146 of (6.8) that this is the probability of the set of states (in Ω) that have
 3147 the property (e.g. ht, th). We say that a random variable x is distributed
 3148 according to a particular probability distribution P_x , which defines the
 3149 probability mapping between the event and the probability of the outcome
 3150 of the random variable. The two concepts are intertwined, but for ease of
 3151 presentation we will discuss some properties with respect to random
 3152 variables and others with respect to their distributions.

3153 *Remark.* The range of the random variable \mathcal{T} is used to indicate the kind
 3154 of probability space, that is a \mathcal{T} random variable. When \mathcal{T} is finite or
 3155 countably infinite, this is called a discrete random variable (Section 6.2.1).
 3156 For continuous random variables (Section 6.2.2) we only consider $\mathcal{T} = \mathbb{R}$
 3157 or $\mathcal{T} = \mathbb{R}^d$. \diamond

3158 6.1.3 Statistics

3159 Probability theory and statistics are often presented together, but they concern
 3160 different aspects of uncertainty. One way of contrasting them is by the
 3161 kinds of problems that are considered. Using probability we can consider
 3162 a model of some process where the underlying uncertainty is captured by
 3163 random variables, and we use the rules of probability to derive what happens.
 3164 Using statistics we observe that something has happened, and try
 3165 to figure out the underlying process that explains the observations. In this
 3166 sense machine learning is close to statistics in its goals, that is to construct
 3167 a model that adequately represents the process that generated the data.
 3168 When the machine learning model is a probabilistic model, we can use the
 3169 rules of probability to calculate the “best fitting” model for some data.

3170 Another aspect of machine learning systems is that we are interested
 3171 in generalization error (see Chapter 8). This means that we are actually
 3172 interested in the performance of our system on instances that we will observe
 3173 in future, which are not identical to the instances that we have seen
 3174 so far. This analysis of future performance relies on probability and statistics,
 3175 most of which is beyond what will be presented in this chapter. The
 3176 interested reader is encouraged to look at the books by Shalev-Shwartz
 3177 and Ben-David (2014); Boucheron et al. (2013). We will see more about
 3178 statistics in Chapter 8.

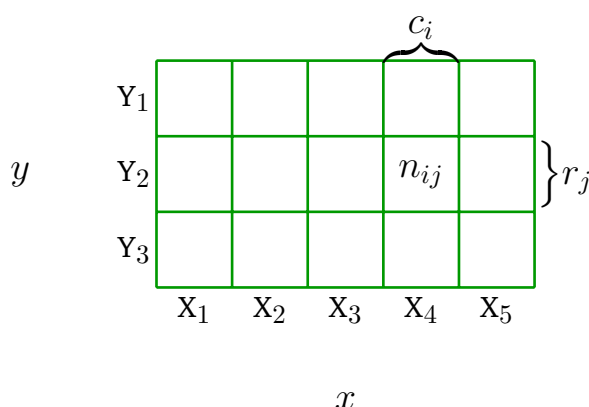


Figure 6.2 Visualization of a discrete bivariate probability mass function, with random variables x and y . This diagram is from Bishop (2006).

6.2 Discrete and Continuous Probabilities

3179

3180 Let us focus our attention on ways to describe the probability of an event
 3181 as introduced in Section 6.1. Depending on whether the state space is
 3182 discrete or continuous the natural way to refer to distributions is different.
 3183 When the outcome space \mathcal{T} is discrete, we can specify the probability that
 3184 a random variable x takes a particular value $x \in \mathcal{T}$, denoted as $P(x = x)$.
 3185 The expression $P(x = x)$ for a discrete random variable x is known as the
 3186 *probability mass function*. When the outcome space \mathcal{T} is continuous, e.g.,
 3187 the real line \mathbb{R} , it is more natural to specify the probability that a random
 3188 variable x is in an interval. By convention we specify the probability that
 3189 a random variable x is less than a particular value x , denoted $P(x \leq x)$.
 3190 The expression $P(x \leq x)$ for a continuous random variable x is known as
 3191 the *cumulative distribution function*. We will discuss continuous random
 3192 variables in Section 6.2.2. We will revisit the nomenclature and contrast
 3193 discrete and continuous random variables in Section 6.2.3.

3194 *Remark.* We will use the phrase *univariate* distribution to refer to distribu-
 3195 tions of only one random variable (denoted by non-bold x). We will refer
 3196 to distributions of more than one random variable as *multivariate* distribu-
 3197 tions, and will usually consider a vector of random variables (denoted
 3198 by bold \boldsymbol{x}). ◇

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

3199

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

joint probability

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

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 $x = x_i$ and $y = Y_j$ is (lazily) written as $p(x, y)$ and is called the joint 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

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

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 $x = \mathbf{x}$, then the fraction of instances (the conditional probability) for which $y = \mathbf{Y}$ 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 labels 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.

categorical variables

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

$$1 \quad \forall \mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) \geq 0$$

2 Its integral exists and

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

3252 Here, $\mathbf{x} \in \mathbb{R}^D$ is a (continuous) random variable. For probability mass
3253 functions (pmf) of discrete random variables the integral in (6.15) is re-
3254 placed with a sum (see (6.12)).

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

3255 In contrast to discrete random variables, the probability of a continuous
3256 random variable x taking a particular value $P(x = x)$ is zero. However
3257 this does not mean that such events never occur.

cumulative
distribution function

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

$$F_{\mathbf{x}}(\mathbf{x}) = P(x_1 \leq x_1, \dots, x_D \leq x_D), \quad (6.16)$$

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

$$F_{\mathbf{x}}(\mathbf{x}) = \int_{-\infty}^{\mathbf{x}} f(\mathbf{x}) d\mathbf{x}. \quad (6.17)$$

3258 6.2.3 Contrasting Discrete and Continuous Distributions

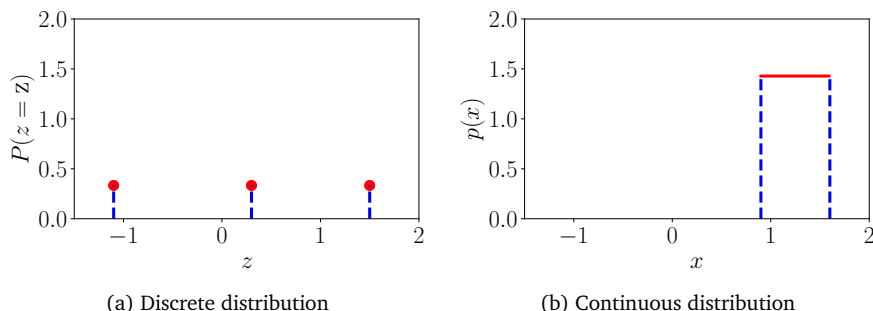
3259 Let us consider both discrete and continuous distributions and contrast
3260 them. The aim here is to see that while both discrete and continuous dis-
3261 tributions seem to have similar requirements, such as the total probability
3262 mass is 1, they are subtly different. Since the total probability mass of a
3263 discrete random variable is 1 (see (6.12)) and there are a finite number of
3264 states, the probability of each state must be in the interval $[0, 1]$. However,
3265 the analogous requirement for continuous random variables (see (6.15))
3266 does not imply that the value of the density is less than or equal to 1 for
uniform distribution 3267 all values. We illustrate this in Figure 6.3 using the *uniform distribution*
3268 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

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



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

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

3269 *Remark.* There is an additional subtlety with regards to discrete prob-
 3270 ability distributions. The states x_1, \dots, x_d do not in principle have any
 3271 structure, that is there is usually no way to compare them, for exam-
 3272 ple $x_1 = \text{red}, x_2 = \text{green}, x_3 = \text{blue}$. However in many machine learn-
 3273 ing applications the discrete states take numerical values, for example
 3274 $x_1 = -1.1, x_2 = 0.3, x_3 = 1.5$, where we could say $x_1 < x_2 < x_3$. Discrete
 3275 states which take numerical values are particularly useful because we of-
 3276 ten consider expected values (Section 6.4.1) of random variables. \diamond

3277 Unfortunately machine learning literature uses notation and nomencla-
 3278 ture that hides the distinction between the state space, the set of outcomes
 3279 and the random variable. For a value x of the set of possible outcomes of
 3280 the random variable x , that is $x \in \mathcal{T}$, $p(x)$ denotes the probability that
 3281 random variable x has the outcome x . For discrete random variables this
 3282 is written as $P(x = x)$, which is known as the probability mass func-

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

tion. 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 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

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} .

sum rule

The first rule, the *sum rule*, is

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

marginalization
property

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

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

3311 by repeated application of the sum rule where we integrate/sum out all
3312 random variables except x_i , which is indicated by $\setminus i$.

3313 *Remark.* Many of the computational challenges of probabilistic modeling
3314 are due to the application of the sum rule. When there are many variables
3315 or discrete variables with many states, the sum rule boils down to per-
3316 forming a high-dimensional sum or integral. Performing high dimensional
3317 sums or integrals is generally computationally hard, in the sense that there
3318 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

product rule

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

3319 The product rule can be interpreted as the fact that every joint distribu-
3320 tion of two random variables can be factorized (written as a product)
3321 of two other distributions. The two factors are the marginal distribu-
3322 tion of the first random variable $p(\mathbf{x})$, and the conditional distribution
3323 of the second random variable given the first $p(\mathbf{y} | \mathbf{x})$. Since the ordering
3324 of random variables is arbitrary in $p(\mathbf{x}, \mathbf{y})$ the product rule also implies
3325 $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x} | \mathbf{y})p(\mathbf{y})$. To be precise, (6.21) is expressed in terms of the
3326 probability mass functions for discrete random variables. For continuous
3327 random variables, the product rule is expressed in terms of the probability
3328 density functions (Section 6.2.3).

3329 Let us briefly explore how to use probabilistic models to capture uncer-
3330 tainty (Ghahramani, 2015). At the lowest modeling level, measurement
3331 noise introduces model uncertainty, for example a camera sensor intro-
3332 duces random error in the value of each pixel it records. We will see in
3333 Chapter 9 how to use Gaussian (Section 6.5) noise models for linear re-
3334 gression. At higher modeling levels, we would be interested in modeling
3335 the uncertainty of the coefficients in linear regression. This uncertainty
3336 captures which values of these parameters will be good at predicting new
3337 data. Finally at the highest levels, we may want to capture uncertainties
3338 about the model structure. We will discuss model choice in Section 8.5.
3339 Once we have the probabilistic models (described in Section 8.3), the
3340 basic rules of probability presented in this section are used to infer the
3341 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

3342 In (6.22), $p(\mathbf{x})$ is the *prior*, which encapsulates our subjective prior
3343 knowledge of the unobserved (latent) variable \mathbf{x} before observing any
3344 data. We can choose any prior that makes sense to us, but it is critical to
3345 ensure that the prior has a non-zero pdf (or pmf) on all plausible \mathbf{x} , even
3346 if they are very rare.

likelihood

The likelihood is sometimes also called the "measurement model".

3347 The *likelihood* $p(\mathbf{y} | \mathbf{x})$ describes how \mathbf{x} and \mathbf{y} are related, and it is the
3348 probability of the data \mathbf{y} if we were to know the latent variable \mathbf{x} . Note
3349 that the likelihood is not a distribution in \mathbf{x} , but only in \mathbf{y} . We call $p(\mathbf{y} | \mathbf{x})$
3350 either the "likelihood of \mathbf{x} (given \mathbf{y})" or the "probability of \mathbf{y} given \mathbf{x} " but
3351 never the likelihood of \mathbf{y} (MacKay, 2003a).

posterior

3352 The *posterior* $p(\mathbf{x} | \mathbf{y})$ is the quantity of interest in Bayesian statistics
3353 because it expresses exactly what we are interested in, i.e., what we know
3354 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
evidence

3355 is the *marginal likelihood/evidence*. By definition the marginal likelihood
3356 integrates the numerator of (6.22) with respect to the latent variable \mathbf{x} .
3357 Therefore, the marginal likelihood is independent of \mathbf{x} and it ensures
3358 that the posterior $p(\mathbf{x} | \mathbf{y})$ is normalized. The marginal likelihood can also
3359 be interpreted as the expected likelihood where we take the expectation
3360 with respect to the prior $p(\mathbf{x})$. Beyond normalization of the posterior the

3361 marginal likelihood also plays an important role in Bayesian model selec-
 3362 tion as we will discuss in Section 8.5. Due to the integration in (8.41), the
 3363 evidence is often hard to compute.

3364 Bayes' theorem in (6.22) allows us to invert the causal relationship be-
 3365 tween x and y given by the likelihood. Therefore, Bayes' theorem is some-
 3366 times called the *probabilistic inverse*.

probabilistic inverse

3367 *Remark.* In Bayesian statistics, the posterior distribution is the quantity
 3368 of interest as it encapsulates all available information from the prior and
 3369 the data. Instead of carrying the posterior around, it is possible to fo-
 3370 cus on some statistic of the posterior, such as the maximum of the pos-
 3371 terior, which we will discuss in Section 9.2.3. However, focusing on the
 3372 some statistic of the posterior leads to loss of information. If we think in
 3373 a bigger context, then the posterior can be used within a decision mak-
 3374 ing system, and having the full posterior around can be extremely useful
 3375 and lead to decisions that are robust to disturbances. For example, in the
 3376 context of model-based reinforcement learning, Deisenroth et al. (2015)
 3377 show that using the full posterior distribution of plausible transition func-
 3378 tions leads to very fast (data/sample efficient) learning, whereas focusing
 3379 on the maximum of the posterior leads to consistent failures. Therefore,
 3380 having the full posterior around in a downstream task can be very use-
 3381 ful. In Chapter 9, we will continue this discussion in the context of linear
 3382 regression. \diamond

3383 6.4 Summary Statistics and Independence

3384 We are often interested in summarizing sets of random variables and com-
 3385 paring pairs of random variables. A statistic of a random variable is a de-
 3386 terministic function of that random variable. The summary statistics of a
 3387 distribution provides one useful view of how a random variable behaves,
 3388 and as the name suggests, provides numbers that summarize and charac-
 3389 terize the distribution. We describe the mean and the variance, two well-
 3390 known summary statistics. Then we discuss two ways to compare a pair
 3391 of random variables: first how to say that two random variables are inde-
 3392 pendent, and second how to compute an inner product between them.

3393 6.4.1 Means and Covariances

3394 Mean and (co)variance are often useful to describe properties of probabili-
 3395 ty distributions (expected values and spread). We will see in Section 6.6
 3396 that there is a useful family of distributions (called the exponential fam-
 3397 ily), where the statistics of the random variable capture all possible infor-
 3398 mation.

3399 The main tool we use to compute statistics of a random variable is its
 3400 expected value with respect to a particular function.

expected value

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

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

where \mathcal{T} is the set of possible outcomes of the random variable x .

In this section, we consider discrete random variables to have outcomes which are numerical. This can be seen by observing that the function g takes real numbers as inputs.

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

The expected value of a function of a random variable is sometimes referred to as the “law of the unconscious statistician” (Casella and Berger, 2002, Section 2.2).
mean

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

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

for $d = 1, \dots, D$, where the subscript d indicates the corresponding dimension of \mathbf{x} , and the integral and sum are over the states x_i .

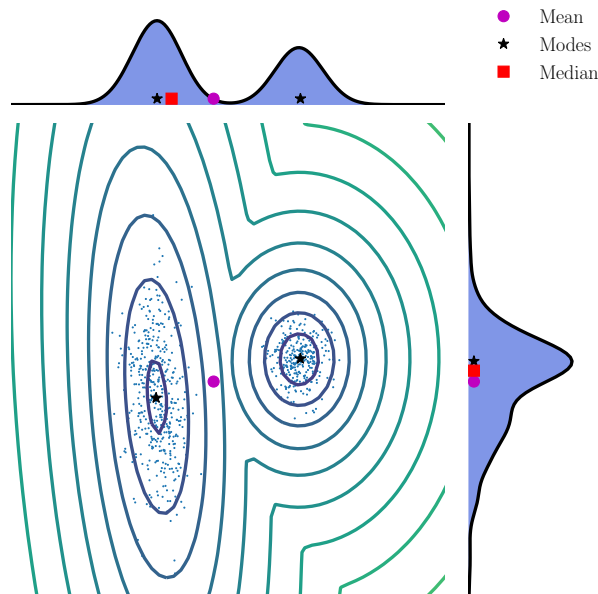


Figure 6.4
Illustration of the mean, mode and median for a two-dimensional dataset, as well as its marginal densities.

3415 In one dimension, there are two other intuitive notions of “average”,
 3416 which are the *median* and the *mode*. The median is the “middle” value if
 3417 we sort the values, i.e., 50% of the values are greater than the median and
 3418 50% are smaller than the median. This idea can be generalized to contin-
 3419 uous values by considering the value where the CDF (Definition 6.2) is
 3420 0.5. For distributions which are asymmetric or have long tails, the median
 3421 provides an estimate of a typical value that is closer to human intuition
 3422 than the mean value. Furthermore the median is more robust to outliers
 3423 than the mean. The generalization of the median to higher dimensions is
 3424 non-trivial as there is no obvious way to “sort” in more than one dimen-
 3425 sion (Hallin et al., 2010; Kong and Mizera, 2012). The *mode* is the most
 3426 frequently occurring value. For a discrete random variable, the mode is
 3427 defined as the value of x having the highest frequency of occurrence. For
 3428 a continuous random random variable, the mode is defined as a peak in
 3429 the density $p(\mathbf{x})$. A particular density $p(\mathbf{x})$ may have more than one mode,
 3430 and, therefore, finding the mode may be computationally challenging in
 3431 high dimensions.

median
mode

mode

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 there is more than one way to define $<$ such that $\begin{bmatrix} 1 \\ 0 \end{bmatrix} < \begin{bmatrix} 2 \\ 3 \end{bmatrix}$.

3432 The mean is recovered if we set the function g in Definition 6.3 to the
 3433 identity function. This indicates that we can think about functions of ran-
 3434 dom variables, which we will revisit in Section 6.7.

Remark. The expected value is a linear operator. For example, given multivariate real-valued functions $f(\mathbf{x}) = ag(\mathbf{x}) + bh(\mathbf{x})$ where $a, b \in \mathbb{R}$,

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

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

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

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

3435

◇

3436 For two random variables, we may wish to characterize their correspon-
 3437 dence to each other. The covariance intuitively represents the notion of
 3438 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

3439

standard deviation

3440

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 $\sigma(x)$.

3441

3442

When we want to compare the covariances between different pairs of

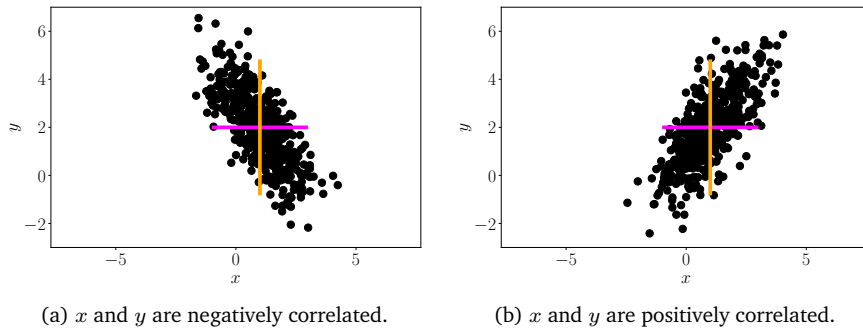


Figure 6.5 Two-dimensional datasets with identical means and variances along each axis (colored lines) but with different covariances.

3443 random variables, it turns out that the variance of each random variable
 3444 affects the value of the covariance. The normalized version of covariance
 3445 is called the correlation.

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

correlation

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

3446 The correlation matrix is the covariance matrix of standardized random
 3447 variables, $x/\sigma(x)$. In other words, each random variable is divided by its
 3448 standard deviation (the square root of the variance) in the correlation
 3449 matrix.

3450 The covariance (and correlation) indicate how two random variables
 3451 are related, see Figure 6.5. Positive correlation $\text{corr}[x, y]$ means that when
 3452 x grows then y is also expected to grow. Negative correlation means that
 3453 as x increases then y decreases.

3454 The notion of covariance can be generalised to multivariate random
 3455 variables.

Terminology:
 Covariance of multivariate random variables $\text{Cov}[x, y]$ is sometimes referred to as cross-covariance, with covariance referring to $\text{Cov}[x, x]$.
 covariance

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

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

3456 Here, the subscript makes it explicit with respect to which variable we
 3457 need to average.

3458 Definition 6.7 can be applied with the same multivariate random vari-
 3459 able in both arguments, which results in a useful concept that intuitively
 3460 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

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

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

covariance matrix 3461

This matrix is called the *covariance matrix* of the random variable \mathbf{x} .

3462

The covariance matrix is symmetric and positive definite and tells us something about the spread of the data.

3463

marginal

The covariance matrix contains the variances of the *marginals*

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

cross-covariance

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

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

3464

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

6.4.2 Empirical Means and Covariances

3465

population mean
and covariance 3466

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.

3472

3473

3474

3475

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.

empirical mean 3477

sample mean 3478

empirical mean

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

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

3479

where $\mathbf{x}_n \in \mathbb{R}^D$.

empirical covariance

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

matrix

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

3480 In the above definition, we have expressed the statistic in terms of the
 3481 multivariate random variables $\mathbf{x}_1, \dots, \mathbf{x}_N$. To compute the statistics for a
 3482 particular dataset, we would use the realizations (observations) X_1, \dots, X_N
 3483 and use (6.41) and (6.42). Empirical covariance matrices are symmetric,
 3484 positive semi-definite (see Section 3.2.3).

3485 *Remark.* The notation we use in this book is imprecise in the sense that
 3486 we do not explicitly make a distinction between the random variable \mathbf{x}_n
 3487 and its deterministic realization x_n . \diamond

Throughout the book we use the empirical covariance, which is a biased estimate. The unbiased (sometimes called corrected) covariance has the factor $N - 1$ in the denominator instead of N .

6.4.3 Three Expressions for the Variance

3488 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.,

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

3489 Depending on whether x is a discrete or continuous random variable, the
 3490 expectation in (6.43) and the mean $\mu = \mathbb{E}_x(x)$ are computed using (6.31).
 3491 The variance as expressed in (6.43) is the mean of a new random variable
 3492 $z := (x - \mu)^2$.

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*:

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

3493 The expression in (6.44) can be remembered as “the mean of the square
 3494 minus the square of the mean”. It can be calculated empirically in one
 3495 pass through data since we can accumulate x_i (to calculate the mean)
 3496 and x_i^2 simultaneously. Unfortunately, if implemented in this way, it can
 3497 be numerically unstable. The raw-score version of the variance can be
 3498 useful in machine learning, e.g., when deriving the bias-variance decom-
 3499 position (Bishop, 2006).

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

raw-score formula for variance

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

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,

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

3500 We see that (6.45) is twice the raw-score expression (6.44). This means
 3501 that we can express the sum of pairwise distances (of which there are N^2
 3502 of them) as a sum of deviations from the mean (of which there are N).
 3503 Geometrically, this means that there is an equivalence between the pair-
 3504 wise distances and the distances from the center of the set of points. From
 3505 a computational perspective, this means that by computing the mean
 3506 (N terms in the summation), and then computing the variance (again
 3507 N terms in the summation) we can obtain an expression (left-hand side
 3508 of (6.45)) that has N^2 terms.

6.4.4 Sums and Transformations of Random Variables

3509
 3510 We may want to model a phenomenon that cannot be well explained by
 3511 textbook distributions (we introduce some in Sections 6.5 and 6.6), and
 3512 hence may perform simple manipulations of random variables (such as
 3513 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

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

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

respectively. Furthermore,

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

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

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

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

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

3514 where $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top$ is the covariance of \mathbf{x} .

3515 6.4.5 Statistical Independence

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

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

3516 Intuitively, two random variables \mathbf{x} and \mathbf{y} are independent if the value
3517 of \mathbf{y} (once known) does not add any additional information about \mathbf{x} (and
3518 vice versa).

3519 If \mathbf{x} , \mathbf{y} are (statistically) independent then

- 3520 • $p(\mathbf{y} | \mathbf{x}) = p(\mathbf{y})$
- 3521 • $p(\mathbf{x} | \mathbf{y}) = p(\mathbf{x})$
- 3522 • $\mathbb{V}_{\mathbf{x}, \mathbf{y}}[\mathbf{x} + \mathbf{y}] = \mathbb{V}_{\mathbf{x}}[\mathbf{x}] + \mathbb{V}_{\mathbf{y}}[\mathbf{y}]$
- 3523 • $\text{Cov}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}, \mathbf{y}] = \mathbf{0}$

3524 The last point above may not hold in converse, that is two random vari-
3525 ables can have covariance zero but are not statistically independent. To
3526 understand why, recall that covariance measures only linear dependence,
3527 therefore random variables that are non-linearly dependent could have
3528 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

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

3529 In machine learning we often consider problems that can be modelled
3530 as *independent and identically distributed* random variables, x_1, \dots, x_N .
3531 The word “independent” refers to Definition 6.10, that is any pair of
3532 random variables x_i and x_j are independent. The phrase *identically dis-*
3533 *tributed* means that all the random variables are from the same distribu-
3534 tion.

independent and
identically
distributed

3535 Another concept that is important in machine learning is conditional
3536 independence.

Definition 6.11 (Conditional Independence). Formally, two random vari-
ables \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}

3537 We write $\mathbf{x} \perp\!\!\!\perp \mathbf{y} \mid z$.

3538 Definition 6.11 requires that the relation in (6.58) must hold true for
 3539 every value of z . The interpretation of (6.58) can be understood as “given
 3540 knowledge about z , the distribution of \mathbf{x} and \mathbf{y} factorizes”. Independence
 3541 can be cast as a special case of conditional independence if we write $\mathbf{x} \perp\!\!\!\perp$
 3542 $\mathbf{y} \mid \emptyset$.

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

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

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

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

3543 Equation (6.60) provides an alternative definition of conditional indepen-
 3544 dence, i.e., $\mathbf{x} \perp\!\!\!\perp \mathbf{y} \mid z$. This alternative presentation provides the interpre-
 3545 tation “given that we know z , knowledge about \mathbf{y} does not change our
 3546 knowledge of \mathbf{x} ”.

3547 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

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

3548 Since variances are measured in squared units, this looks very much like
 3549 the Pythagorean theorem for right triangles $c^2 = a^2 + b^2$.

Inner products
between
multivariate random
variables \mathbf{x}, \mathbf{y} can
be treated in a
similar fashion

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

$\text{Cov}[x, x] > 0$ and
 $0 \iff x = 0$
 $\text{Cov}[\alpha x + z, y] =$
 $\alpha \text{Cov}[x, y] +$
 $\text{Cov}[z, y]$ for $\alpha \in \mathbb{R}$.

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

3550 i.e., its standard deviation. The “longer” the random variable, the more
 3551 uncertain it is; and a random variable with length 0 is deterministic.

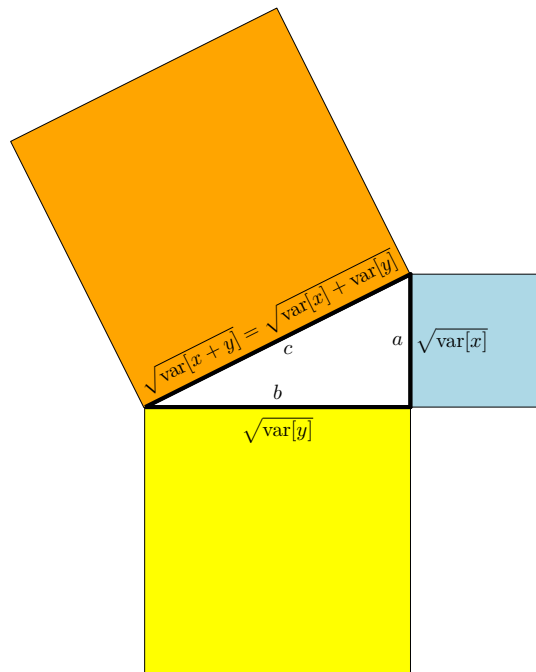


Figure 6.6
 Geometry of random variables. If random variables x and y are uncorrelated they are orthogonal vectors in a corresponding vector space, and the Pythagorean theorem applies.

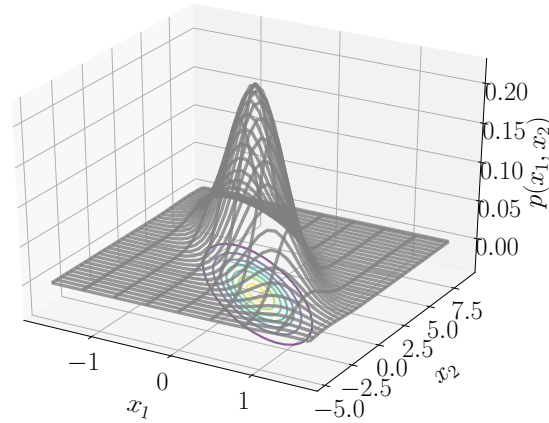
If we look at the angle θ between random two random variables x, y , we get

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

3552 which is the correlation (Definition 6.6) between the two random vari-
 3553 ables. This means that we can think of correlation as the angle between
 3554 two random variables when we consider them geometrically. We know
 3555 from Definition 3.7 that $x \perp y \iff \langle x, y \rangle = 0$. In our case this means
 3556 that x and y are orthogonal if and only if $\text{Cov}[x, y] = 0$, i.e., they are
 3557 uncorrelated. Figure 6.6 illustrates this relationship.

3558 *Remark.* While it is tempting to use the Euclidean distance (constructed
 3559 from the definition of inner products above) to compare probability distri-
 3560 butions, it is unfortunately not the best way to obtain distances between
 3561 distributions. Recall that the probability mass (or density) is positive and
 3562 needs to add up to 1. These constraints mean that distributions live on
 3563 something called a statistical manifold. The study of this space of prob-
 3564 ability distributions is called information geometry. Computing distances
 3565 between distributions are often done using Kullback-Leibler divergence
 3566 which is a generalization of distances that account for properties of the
 3567 statistical manifold. Just like the Euclidean distance is a special case of a
 3568 metric (Section 3.3) the Kullback-Leibler divergence is a special case of
 3569 two more general classes of divergences called Bregman divergences and

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



3570 f -divergences. The study of divergences is beyond the scope of this book.
 3571 Interested readers are referred to a recent book (Amari, 2016) written by
 3572 one of the founders of the field of information geometry. \diamond

6.5 Gaussian Distribution

The Gaussian 3573
 distribution arises 3574
 naturally when we 3575
 consider sums of 3576
 independent and 3577
 identically 3578
 distributed random 3579
 variables. This is 3580
 known as the 3581
 Central Limit 3582
 Theorem (Grinstead 3583
 and Snell, 1997). 3584
 normal distribution 3585

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 density that is given by

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

multivariate 3586
 Gaussian 3587
 distribution
 Also: multivariate
 normal distribution
 mean vector
 covariance matrix

The *multivariate Gaussian distribution* is fully characterized by a *mean vector* μ and a *covariance matrix* Σ and defined as

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

where $\mathbf{x} \in \mathbb{R}^D$ is a random variable. We write $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \mu, \Sigma)$ or $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$. Figure 6.7 shows a bi-variate Gaussian (mesh), with the

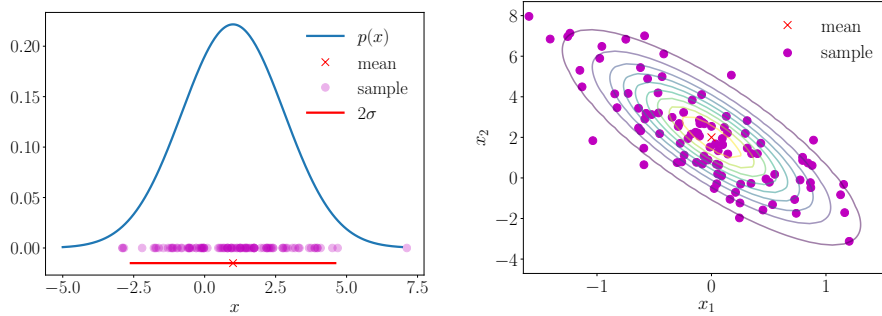


Figure 6.8 Gaussian distributions overlaid with 100 samples. Left: Univariate (1-dimensional) Gaussian; The red cross shows the mean and the red line shows the extent of the variance. Right: Multivariate (2-dimensional) Gaussian, viewed from top. The red standard normal cross shows the distribution mean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

3588 corresponding contour plot. The special case of the Gaussian with zero
 3589 mean and identity covariance, that is $\mu = \mathbf{0}$ and $\Sigma = \mathbf{I}$, is referred to as
 3590 the *standard normal distribution*.

3591 Gaussian distributions are widely used in statistical estimation and machine
 3592 learning because they have closed-form expressions for marginal
 3593 and conditional distributions. In Chapter 9, we use these closed form
 3594 expressions extensively for linear regression. A major advantage of modelling
 3595 with Gaussian distributed random variables is that variable transformations
 3596 (Section 6.7) are often not needed. Since the Gaussian distribution is fully
 3597 specified by its mean and covariance we often can obtain the transformed
 3598 distribution by applying the transformation to the mean and covariance of the
 3599 random variable.

3600 **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 $[\mathbf{x}, \mathbf{y}]^\top$,

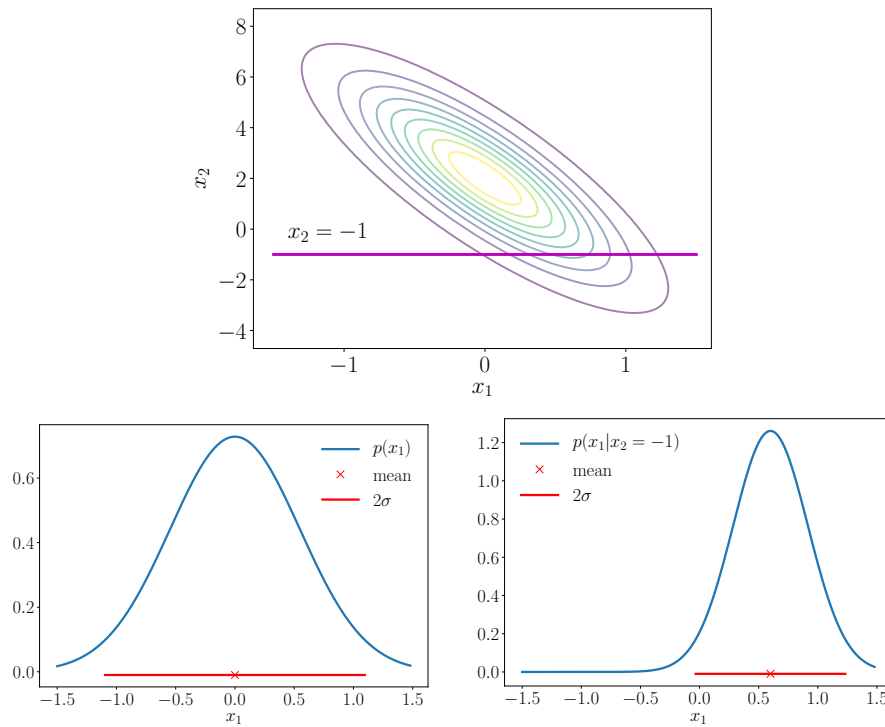
$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right). \tag{6.67}$$

3601 where $\Sigma_{xx} = \text{Cov}[\mathbf{x}, \mathbf{x}]$ and $\Sigma_{yy} = \text{Cov}[\mathbf{y}, \mathbf{y}]$ are the marginal covariance
 3602 matrices of \mathbf{x} and \mathbf{y} , respectively, and $\Sigma_{xy} = \text{Cov}[\mathbf{x}, \mathbf{y}]$ is the cross-
 3603 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))

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mu_{x|y}, \Sigma_{x|y}) \tag{6.68}$$

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



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

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

3604 Note that in the computation of the mean in (6.69) the \mathbf{y} -value is an
3605 observation and no longer random.

3606 *Remark.* The conditional Gaussian distribution shows up in many places,
3607 where we are interested in posterior distributions:

- 3608 • The Kalman filter (Kalman, 1960), one of the most central algorithms
3609 for state estimation in signal processing, does nothing but computing
3610 Gaussian conditionals of joint distributions (Deisenroth and Ohlsson,
3611 2011).
- 3612 • Gaussian processes (Rasmussen and Williams, 2006), which are a practical
3613 implementation of a distribution over functions. In a Gaussian process,
3614 we make assumptions of joint Gaussianity of random variables. By
3615 (Gaussian) conditioning on observed data, we can determine a posterior
3616 distribution over functions.
- 3617 • Latent linear Gaussian models (Roweis and Ghahramani, 1999; Murphy,
3618 2012), which include probabilistic PCA (Tipping and Bishop, 1999).

3619

◇

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

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

3620 The corresponding result holds for $p(\mathbf{y})$, which is obtained by marginaliz-
3621 ing with respect to \mathbf{x} . Intuitively, looking at the joint distribution in (6.67),
3622 we ignore (i.e., integrate out) everything we are not interested in. This is
3623 illustrated on the bottom left of Figure 6.9.

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

3624 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

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

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

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

3625 The scaling constant c itself can be written in the form of a Gaussian
 3626 density either in \mathbf{a} or in \mathbf{b} with an “inflated” covariance matrix $\mathbf{A} + \mathbf{B}$,
 3627 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

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

3628 Here, neither \mathbf{a} nor \mathbf{b} are random variables. However, writing c in this way
 3629 is more compact than (6.79). \diamond

3630 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

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

3631 Knowing that $p(\mathbf{x} + \mathbf{y})$ is Gaussian, the mean and covariance matrix can be
 3632 determined immediately using the results from (6.46)–(6.49). This prop-
 3633 erty will be important when we consider i.i.d. Gaussian noise acting on
 3634 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

$$p(a\mathbf{x} + b\mathbf{y}) = \mathcal{N}(a\boldsymbol{\mu}_x + b\boldsymbol{\mu}_y, a^2\boldsymbol{\Sigma}_x + b^2\boldsymbol{\Sigma}_y). \quad (6.82)$$

3635 *Remark.* A case which will be useful in Chapter 11 is the weighted sum
 3636 of Gaussian densities. This is different from the weighted sum of Gaussian
 3637 random variables. \diamond

3638 In Theorem 6.12, the random variable x is from a density which a mix-
 3639 ture of two densities $p_1(x)$ and $p_2(x)$, weighted by α . The theorem can
 3640 be generalized to the multivariate random variable case, since linearity of
 3641 expectations holds also for multivariate random variables. However the
 3642 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 densi-
 ties

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

3643 where the scalar $0 < \alpha < 1$ is the mixture weight, and $p_1(x)$ and $p_2(x)$ are

3644 univariate Gaussian densities (Equation (6.65)) with different parameters,
3645 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).

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

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

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

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

3646 where in the last equality, we again used the raw score version of the vari-
3647 ance and rearranged terms such that the expectation of a squared random
3648 variable is the sum of the squared mean and the variance.

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

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

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

$$= [\alpha\sigma_1^2 + (1 - \alpha)\sigma_2^2]$$

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

3649 For a mixture, the individual components can be considered to be condi-
 3650 tional distributions (conditioned on the component identity). The last line
 3651 is an illustration of the conditional variance formula: “The variance of a
 3652 mixture is the mean of the conditional variance and the variance of the
 3653 conditional mean”. \square

3654 *Remark.* The derivation above holds for any density, but in the case of
 3655 the Gaussian since it is fully determined by the mean and variance, the
 3656 mixture density can be determined in closed form. \diamond

3657 We consider in Example 6.18 a bivariate standard Gaussian random
 3658 variable \mathbf{x} and performed a linear transformation $\mathbf{A}\mathbf{x}$ on it. The outcome
 3659 is a Gaussian random variable with zero mean and covariance $\mathbf{A}\mathbf{A}^\top$. Ob-
 3660 serve that adding a constant vector will change the mean of the distribu-
 3661 tion, without affecting its variance, that is the random variable $\mathbf{x} + \boldsymbol{\mu}$ is
 3662 Gaussian with mean $\boldsymbol{\mu}$ and identity covariance. Therefore, a linear (or
 3663 affine) transformation of a Gaussian random variable is Gaussian dis-
 3664 tributed.

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

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

3665 **6.5.4 Sampling from Multivariate Gaussian Distributions**

3666 We will not explain the subtleties of random sampling on a computer. In
 3667 the case of a multivariate Gaussian, this process consists of three stages:
 3668 first we need a source of pseudo-random numbers that provide a uniform
 3669 sample in the interval $[0,1]$, second we use a non-linear transformation
 3670 such as the Box-Müller transform (Devroye, 1986) to obtain a sample from
 3671 a univariate Gaussian, and third we collate a vector of these samples to
 3672 obtain a sample from a multivariate standard normal $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

3673 For a general multivariate Gaussian, that is where the mean is non-zero
 3674 and the covariance is not the identity matrix, we use the properties of
 3675 linear transformations of a Gaussian random variable. Assume we are inter-
 3676 ested in generating samples $\mathbf{x}_i, i = 1, \dots, n$, from a multivariate Gaus-
 3677 sian distribution with mean $\boldsymbol{\mu}$ and covariance matrix Σ . We would like
 3678 to construct the sample from a sampler that provides samples from the
 3679 multivariate standard normal $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

3680 To obtain samples from a multivariate normal $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$, we can use
 3681 the properties of a linear transformation of a Gaussian random variable:
 3682 If $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ then $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\mu}$, where $\mathbf{A}\mathbf{A}^\top = \Sigma$, is Gaussian dis-
 3683 tributed with mean $\boldsymbol{\mu}$ and covariance matrix Σ . Recall from Section 4.3
 3684 that when we can decompose $\Sigma = \mathbf{A}\mathbf{A}^\top$, while there are many possible
 3685 decompositions, we often choose the Cholesky decomposition. This has
 3686 the benefit that \mathbf{A} is triangular, leading to efficient computation.

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.

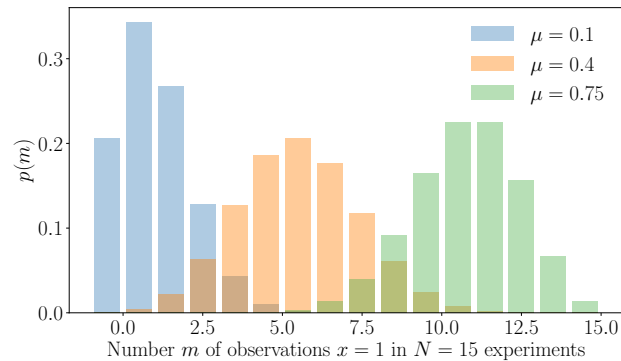
3687 **6.6 Conjugacy and the Exponential Family**

3688 Many of the probability distributions “with names” that we find in statis-
 3689 tics textbooks were discovered to model particular types of phenomena.
 3690 For example we have seen the Gaussian distribution in Section 6.5. The
 3691 distributions are also related to each other in complex ways (Leemis and
 3692 McQueston, 2008). For a beginner in the field, it can be overwhelming to
 3693 figure out which distribution to use. In addition, many of these distribu-
 3694 tions were discovered at a time that statistics and computation was done
 3695 by pencil and paper. It is natural to ask what are meaningful concepts
 3696 in the computing age (Efron and Hastie, 2016). In the previous section,
 3697 we saw that many of the operations required for inference can be conven-
 3698 iently calculated when the distribution is Gaussian. It is worth recalling
 3699 at this point the desiderata for manipulating probability distributions.

“Computers” were a job description.

3700 1 There is some “closure property” when applying the rules of probability,
 3701 e.g., Bayes’ theorem. By closure we mean that applying a particular
 3702 operation returns an object of the same type.

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



- 3703 2 As we collect more data, we do not need more parameters to describe
3704 the distribution.
3705 3 Since we are interested in learning from data, we want parameter esti-
3706 mation to behave nicely.

exponential family 3707 It turns out that the class of distributions called the *exponential family*
3708 provides the right balance of generality while retaining favourable com-
3709 putation and inference properties. Before we introduce the exponential
3710 family, let us see three more members of “named” probability distribu-
3711 tions, the Bernoulli (Example 6.8), Binomial (Example 6.9) and Beta (Ex-
3712 ample 6.10) distributions.

Bernoulli
distribution



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

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

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

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

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

3713 An example where the Bernoulli distribution can be used is when we
3714 are interested in modeling the probability of “head” when flipping a coin.

3715 *Remark.* The rewriting above of the Bernoulli distribution, where we use
3716 Boolean variables as numerical 0 or 1 and express them in the exponents,
3717 is a trick that is often used in machine learning textbooks. Another oc-
3718 currence of this is when expressing the Multinomial distribution. \diamond

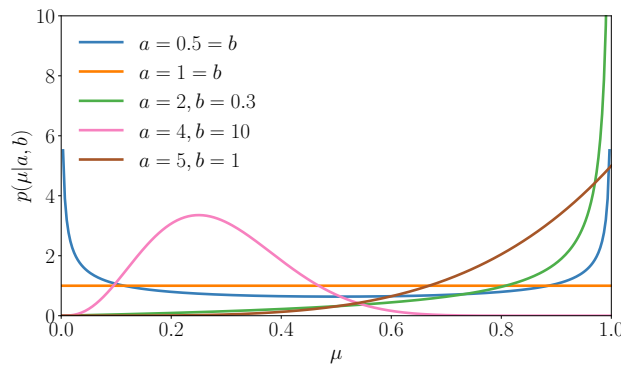


Figure 6.11
Examples of the Beta distribution for different values of α and β .

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

Binomial distribution

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

$$\mathbb{E}[m] = N\mu, \tag{6.99}$$

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

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

3719 An example where the Binomial could be used is if we want to describe
3720 the probability of observing m “heads” in N coin-flip experiments if the
3721 probability for observing head in a single experiment is μ .

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

Beta distribution

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

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

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.

3722 Intuitively, α moves probability mass toward 1, whereas β moves prob-
3723 ability mass toward 0. There are some special cases (Murphy, 2012):

- 3724 • For $\alpha = 1 = \beta$ we obtain the uniform distribution $\mathcal{U}[0, 1]$.
- 3725 • For $\alpha, \beta < 1$, we get a bimodal distribution with spikes at 0 and 1.
- 3726 • For $\alpha, \beta > 1$, the distribution is unimodal.
- 3727 • For $\alpha, \beta > 1$ and $\alpha = \beta$, the distribution is unimodal, symmetric and
3728 centered in the interval $[0, 1]$, i.e., the mode/mean is at $\frac{1}{2}$.

3729 *Remark.* There is a whole zoo of distributions with names, and they are
3730 related in different ways to each other (Leemis and McQueston, 2008).
3731 It is worth keeping in mind that each named distribution is created for a
3732 particular reason, but may have other applications. Knowing the reason
3733 behind the creation of a particular distribution often allows insight into
3734 how to best use it. We introduced the above three distributions to be able
3735 to illustrate the concepts of conjugacy (Section 6.6.1) and exponential
3736 families (Section 6.6.3). \diamond

3737 6.6.1 Conjugacy

3738 According to Bayes' theorem (6.22), the posterior is proportional to the
3739 product of the prior and the likelihood. The specification of the prior can
3740 be tricky for two reasons: First, the prior should encapsulate our knowl-
3741 edge about the problem before we see some data. This is often difficult to
3742 describe. Second, it is often not possible to compute the posterior distribu-
3743 tion analytically. However, there are some priors that are computationally
3744 convenient: *conjugate priors*.

conjugate 3745 **Definition 6.13** (Conjugate Prior). A prior is *conjugate* for the likelihood
3746 function if the posterior is of the same form/type as the prior.

3747 Conjugacy is particularly convenient because we can algebraically cal-
3748 culate our posterior distribution by updating the parameters of the prior
3749 distribution.

3750 *Remark.* When considering the geometry of probability distributions, con-
3751 jugate priors retain the same distance structure as the likelihood (Agarwal
3752 and III, 2010). \diamond

3753 To introduce a concrete example of conjugate priors, we describe below
 3754 the Binomial distribution (defined on discrete random variables) and the
 3755 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

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

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

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

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

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

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

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

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

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

Table 6.2 Examples of conjugate priors for common likelihood functions.

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

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 (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).

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.

6.6.2 Sufficient Statistics

Recall that a statistic of a random variable is a deterministic function of that random variable. For example if $\mathbf{x} = [x_1, \dots, x_N]^T$ 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 $\phi(x)$.

3787 The Fisher-Neyman factorization theorem formalizes this notion, which
3788 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*

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

3789 where $h(x)$ is a distribution independent of θ and g_{θ} captures all the depen-
3790 dence on θ via sufficient statistics $\phi(x)$.

3791 If $p(x | \theta)$ does not depend on θ then $\phi(x)$ is trivially a sufficient statistic
3792 for any function ϕ . The more interesting case is that $p(x | \theta)$ is dependent
3793 only on $\phi(x)$ and not x itself. In this case, $\phi(x)$ is a sufficient statistic for
3794 x .

3795 In machine learning we consider a finite number of samples from a
3796 distribution. One could imagine that for simple distributions (such as the
3797 Bernoulli in Example 6.8) we only need a small number of samples to
3798 estimate the parameters of the distributions. We could also consider the
3799 opposite problem: if we have a set of data (a sample from an unknown
3800 distribution), which distribution gives the best fit? A natural question to
3801 ask is as we observe more data, do we need more parameters θ to de-
3802 scribe the distribution? It turns out that the answer is yes in general, and
3803 this is studied in non-parametric statistics (Wasserman, 2007). A converse
3804 question is to consider which class of distributions have finite dimensional
3805 sufficient statistics, that is the number of parameters needed to describe
3806 them do not increase arbitrarily. The answer is exponential family distri-
3807 butions, described in the following section.

3808 6.6.3 Exponential Family

3809 There are three possible levels of abstraction we can have when con-
3810 sidering distributions (of discrete or continuous random variables). At
3811 level one (the most concrete end of the spectrum), we have a particu-
3812 lar named distribution with fixed parameters, for example a univariate
3813 Gaussian $\mathcal{N}(0, 1)$ with zero mean and unit variance. In machine learning
3814 we often use the second level of abstraction, that is we fix the paramet-
3815 ric form (the univariate Gaussian) and infer the parameters from data. For
3816 example, we assume a univariate Gaussian $\mathcal{N}(\mu, \sigma^2)$ with unknown mean
3817 μ and unknown variance σ^2 , and use a maximum likelihood fit to deter-
3818 mine the best parameters (μ, σ^2) . We will see an example of this when
3819 considering linear regression in Chapter 9. A third level of abstraction is
3820 to consider families of distributions, and in this book, we consider the ex-
3821ponential family. The univariate Gaussian is an example of a member of
3822 the exponential family. Many of the widely used statistical models, includ-

3823 ing all the “named” models in Table 6.2, are members of the exponential
3824 family. They can all be unified into one concept (Brown, 1986).

3825 *Remark.* A brief historical anecdote: like many concepts in mathematics
3826 and science, exponential families were independently discovered at the
3827 same time by different researchers. In the years 1935–1936, Edwin Pitman
3828 in Tasmania, Georges Darmon in Paris, and Bernard Koopman in New
3829 York, independently showed that the exponential families are the only
3830 families that enjoy finite-dimensional sufficient statistics under repeated
3831 independent sampling (Lehmann and Casella, 1998). \diamond

exponential family

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

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

3832 where $\boldsymbol{\phi}(\boldsymbol{x})$ is the vector of sufficient statistics. In general, any inner prod-
3833 uct (Section 3.2) can be used in (6.110), and for concreteness we will use
3834 the standard dot product here ($\langle \boldsymbol{\theta}, \boldsymbol{\phi}(\boldsymbol{x}) \rangle = \boldsymbol{\theta}^\top \boldsymbol{\phi}(\boldsymbol{x})$). Note that the form
3835 of the exponential family is essentially a particular expression of $g_\theta(\boldsymbol{\phi}(\boldsymbol{x}))$
3836 in the Fisher-Neyman theorem (Theorem 6.14).

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

log partition function

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

natural parameters

3837 For this form of parameterization, the parameters $\boldsymbol{\theta}$ are called the *natural*
3838 *parameters*. At first glance it seems that exponential families is a mundane
3839 transformation by adding the exponential function to the result of a dot
3840 product. However, there are many implications that allow for convenient
3841 modelling and efficient computation based on the fact that we can capture
3842 information about data in $\boldsymbol{\phi}(\boldsymbol{x})$.

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 $\phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$.

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).

3843 *Remark.* The relationship between the original Bernoulli parameter μ and
 3844 the natural parameter θ is known as the *sigmoid* or logistic function. Ob-
 3845 serve that $\mu \in (0, 1)$ but $\theta \in \mathbb{R}$, and therefore the sigmoid function
 3846 squeezes a real value into the range $(0, 1)$. This property is useful in ma-
 3847 chine learning, for example it is used in logistic regression (Bishop, 1995,
 3848 Section 4.3.2), as well as as a nonlinear activation functions in neural
 3849 networks (Goodfellow et al., 2016, Chapter 6). \diamond sigmoid

It is often not obvious how to find the parametric form of the conju-

gate distribution of a particular distribution. Exponential families provide a convenient way to find conjugate pairs of distributions. Consider the random variable \mathbf{x} distributed as an exponential family (6.110)

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

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

$$p(\boldsymbol{\theta} | \boldsymbol{\gamma}) = h_c(\boldsymbol{\theta}) \exp\left(\left\langle \begin{bmatrix} \boldsymbol{\gamma}_1 \\ \boldsymbol{\gamma}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\theta} \\ -A(\boldsymbol{\theta}) \end{bmatrix} \right\rangle - A_c(\boldsymbol{\gamma})\right) , \quad (6.126)$$

3850 where $\boldsymbol{\gamma} = \begin{bmatrix} \boldsymbol{\gamma}_1 \\ \boldsymbol{\gamma}_2 \end{bmatrix}$ has dimension $\dim(\boldsymbol{\theta}) + 1$. The sufficient statistics of
 3851 the conjugate prior are $\begin{bmatrix} \boldsymbol{\theta} \\ -A(\boldsymbol{\theta}) \end{bmatrix}$. By using the knowledge of the general
 3852 form of conjugate priors for exponential families, we can derive functional
 3853 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.

3854 As mentioned in the previous section, the main motivation for expo-
 3855 nential families is that they have finite-dimensional sufficient statistics.
 3856 Additionally, conjugate distributions are easy to write down, and the con-
 3857 jugate distributions also come from an exponential family. From an infer-
 3858 ence perspective, maximum likelihood estimation behaves nicely because
 3859 empirical estimates of sufficient statistics are optimal estimates of the pop-
 3860 ulation values of sufficient statistics (recall the mean and covariance of a

3861 Gaussian). From an optimization perspective, the log-likelihood function
 3862 is concave allowing for efficient optimization approaches to be applied
 3863 (Chapter 7).

3864 6.7 Change of Variables/Inverse Transform

3865 It may seem that there are very many known distributions, but in re-
 3866 ality the set of distributions for which we have names is quite limited.
 3867 Therefore, it is often useful to understand how transformed random vari-
 3868 ables are distributed. For example, assume that x is a random variable
 3869 distributed according to the univariate normal distribution $\mathcal{N}(0, 1)$, what
 3870 is the distribution of x^2 ? Another example, which is quite common in ma-
 3871 chine learning, is: given that x_1 and x_2 are univariate standard normal,
 3872 what is the distribution of $\frac{1}{2}(x_1 + x_2)$?

3873 One option to work out the distribution of $\frac{1}{2}(x_1 + x_2)$ is to calculate the
 3874 mean and variance of x_1 and x_2 and then combine them. As we saw in
 3875 Section 6.4.4, we can calculate the mean and variance of resulting ran-
 3876 dom variables when we consider affine transformations of random vari-
 3877 ables. However, we may not be able to obtain the functional form of the
 3878 distribution under transformations. Furthermore, we may be interested
 3879 in nonlinear transformations of random variables for which closed-form
 3880 expressions are not readily available.

3881 *Remark (Notation).* In this section, we will be explicit about random vari-
 3882 ables and the values they take. Hence, we will use small letters x, y to
 3883 denote random variables and small capital letters x, Y to denote the values
 3884 that the random variables take. We will explicitly write probability mass
 3885 functions (pmf) of discrete random variables x as $P(x = X)$. For contin-
 3886 uous random variables x , the probability density function (pdf) is written
 3887 as $f(x)$ and the the cumulative distribution function (cdf) is written as
 3888 $F_x(x \leq X)$. \diamond

3889 We will look at two approaches for obtaining distributions of transfor-
 3890 mations of random variables: a direct approach using the definition of a
 3891 cumulative distribution function and a change-of-variable approach that
 3892 uses the chain rule of calculus (Section 5.2.2). The change-of-variable ap-
 3893 proach is widely used because it provides a “recipe” for attempting to
 3894 compute the resulting distribution due to a transformation. We will ex-
 3895 plain the techniques for univariate random variables, and will only briefly
 3896 provide the results for the general case of multivariate random variables.

3897 As mentioned in Section 6.1, random variables and probability distri-
 3898 butions are closely associated with each other. It is worth carefully teasing
 3899 apart the two ideas, and in doing so we will motivate why we need to
 3900 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

$$\begin{aligned} P(y = Y) &= P(U(x) = Y) && \text{transformation of interest} && (6.131a) \\ &= P(x = U^{-1}(Y)) && \text{inverse} && (6.131b) \end{aligned}$$

3901 where we can observe that $x = U^{-1}(Y)$. Therefore for discrete random
3902 variables, transformations directly change the individual events (with the
3903 probabilities appropriately transformed).

6.7.1 Distribution Function Technique

3904
3905 The distribution function technique goes back to first principles, and uses
3906 the definition of a cumulative distribution function (cdf) $F_x(x) = P(x \leq$
3907 $x)$ and the fact that its differential is the probability density function (pdf)
3908 $f(x)$ (Wasserman, 2004, Chapter 2). For a random variable x and a func-
3909 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 $f(y)$.

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

3910 We also need to keep in mind that the domain of the random variable may
3911 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 $0 \leq y \leq 1$.

3912 In Example 6.17, we considered a strictly monotonically increasing function
 3913 $f(x) = 3x^2$. This means that we could compute an inverse function.
 3914 In general, we require that the function of interest $y = U(x)$ has an inverse
 3915 $x = U^{-1}(y)$. A useful result can be obtained by considering the
 3916 cumulative distribution function $F_x(x)$ of a random variable x , and using
 3917 it as the transformation $U(x)$. This leads to the following theorem
 3918 which is called the probability integral transform. The result is the basis
 3919 of generating samples from distributions whose cdfs are known by first
 3920 generating a sample from a uniform distribution and then transforming it
 3921 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*

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

3922 *has a uniform distribution.*

3923 *Proof* We need to show that the cumulative distribution function of y
 3924 defines a distribution of a uniform random variable. Recall that by the
 3925 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

probability integral
transform

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).

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.

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

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

3953 are somehow “inverses” of each other. An intuitive understanding of the
 3954 rule can be obtained by thinking (loosely) about small changes (differen-
 3955 tials) to the equation $u = g(x)$. That is by considering $\Delta u = g'(x)\Delta x$ as a
 3956 differential of $u = g(x)$. By substituting $u = g(x)$, the argument inside the
 3957 integral on the right hand side of (6.140) becomes $f(g(x))$. By pretending
 3958 that the term du can be approximated by $du \approx \Delta u = g'(x)\Delta x$, and that
 3959 $dx \approx \Delta x$, 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

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

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

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

change of variable³⁹⁶⁰ This is called the *change of variable* technique. The term $\left| \frac{d}{dy} U^{-1}(y) \right|$
³⁹⁶¹ measures how much a unit volume changes when applying U (see also
³⁹⁶² the Remark on page 151. In (6.149) we introduced the absolute value of
³⁹⁶³ the differential. For decreasing functions, it turns out that an additional
³⁹⁶⁴ negative sign is needed, and instead of having two types of change-of-
³⁹⁶⁵ variable rules, the absolute value unifies both of them.

³⁹⁶⁶ *Remark.* In comparison to the discrete case in (6.131b), we have an addi-
³⁹⁶⁷ tional factor $\left| \frac{d}{dy} U^{-1}(y) \right|$. The continuous case requires more care because
³⁹⁶⁸ $P(y = Y) = 0$ for all Y . The probability density function $f(y)$ does not
³⁹⁶⁹ have a description as a probability of an event involving y . \diamond

³⁹⁷⁰ So far in this section we have been studying univariate change of vari-
³⁹⁷¹ ables. The case for multivariate random variables is analogous, but com-
³⁹⁷² plicated by fact that the absolute value cannot be used for multivariate
³⁹⁷³ functions. Instead we use the determinant of the Jacobian matrix. Recall
³⁹⁷⁴ from (5.56) that the Jacobian is a matrix of partial derivatives, and that
³⁹⁷⁵ the existence of a non-zero determinant shows that we can invert the Ja-
³⁹⁷⁶ cobian. Recall the discussion in Section 4.1 that the determinant arises
³⁹⁷⁷ because our differentials (cubes of volume) are transformed into paral-
³⁹⁷⁸ lelepipeds by the Jacobian. Let us summarize the discussion above in the
³⁹⁷⁹ following theorem, which gives us a recipe for multivariate change of vari-
³⁹⁸⁰ 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

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

³⁹⁸¹ The theorem looks intimidating at first glance, but the key point is that
³⁹⁸² a change of variable of a multivariate random variable follows the pro-
³⁹⁸³ cedure of the univariate change of variable. First we need to work out
³⁹⁸⁴ the inverse transform, and substitute that into the density of \mathbf{x} . Then we
³⁹⁸⁵ calculate the determinant of the Jacobian and multiply the result. The
³⁹⁸⁶ 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

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

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

3987 While Example 6.18 is based on a bivariate random variable, which
 3988 allows to easily compute the matrix inverse, the relation above holds for
 3989 higher dimensions.

3990 *Remark.* We saw in Section 6.5 that the density $f(\boldsymbol{x})$ above is actually
 3991 the standard Gaussian distribution, and the transformed density $f(\boldsymbol{y})$ is a
 3992 bivariate Gaussian with covariance $\boldsymbol{\Sigma} = \boldsymbol{A}\boldsymbol{A}^\top$. \diamond

3993 6.8 Further Reading

3994 This chapter is rather terse at times, Grinstead and Snell (1997) and
 3995 Walpole et al. (2011) provides more relaxed presentations that are suit-
 3996 able for self study. Readers interested in more philosophical aspects of
 3997 probability should consider Hacking (2001), whereas a more software
 3998 engineering approach is presented by Downey (2014). An overview of
 3999 exponential families can be found in Barndorff-Nielsen (2014). We will
 4000 see more about how to use probability distributions to model machine
 4001 learning tasks in Chapter 8. Ironically the recent surge in interest in neural
 4002 networks has resulted in a broader appreciation of probabilistic mod-
 4003 els. For example the idea of normalizing flows (Rezende and Mohamed,
 4004 2015) relies on change of variables for transforming random variables.
 4005 An overview of methods for variational inference as applied to neural net-
 4006 works is described in Chapters 16 to 20 of Goodfellow et al. (2016).

4007 We side stepped a large part of the difficulty in continuous random vari-
 4008 ables by avoiding measure theoretic questions (Billingsley, 1995; Pollard,
 4009 2002), and by assuming without construction that we have real numbers,
 4010 and ways of defining sets on real numbers as well as their appropriate fre-
 4011 quency of occurrence. These details do matter, for example in the specifi-
 4012 cation of conditional probability $p(y|x)$ for continuous random variables
 4013 x, y (Proschan and Presnell, 1998). The lazy notation hides the fact that
 4014 we want to specify that $x = \mathbf{x}$ (which is a set of measure zero). Further-
 4015 more we are interested in the probability density function of y . A more
 4016 precise notation would have to say $\mathbb{E}_y[f(y) | \sigma(x)]$, where we take the
 4017 expectation over y of a test function f conditioned on the σ -algebra of
 4018 x . A more technical audience interested in the details of probability the-
 4019 ory have many options (Jacod and Protter, 2004; Jaynes, 2003; MacKay,
 4020 2003b; Grimmett and Welsh, 2014) including some very technical dis-
 4021 cussions (Çinlar, 2011; Dudley, 2002; Shiryaev, 1984; Lehmann and
 4022 Casella, 1998; Bickel and Doksum, 2006). As machine learning allows
 4023 us to model more intricate distributions on ever more complex types of
 4024 data, a developer of probabilistic machine learning models would have to
 4025 understand these more technical aspects. Machine learning books with a
 4026 probabilistic modelling focus includes MacKay (2003b); Bishop (2006);
 4027 Murphy (2012); Barber (2012); Rasmussen and Williams (2006).

Exercises

4028 6.1 Consider a mixture of two Gaussian distributions (illustrated in Figure 6.4)

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

- 4029 1 Compute the marginal distributions for each dimension
 4030 2 Compute the mean, mode and median for each marginal distribution
 4031 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 :

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

4032 Choose a conjugate prior for the Bernoulli likelihood and compute the posterior distribution $p(\mu|x_1, \dots, x_N)$.

4033 6.3 Consider the following time-series model:

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

4034 where \mathbf{w}, \mathbf{v} are i.i.d. Gaussian noise variables. Further, assume that $p(\mathbf{x}_0) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$.

- 4036 1 What is the form of $p(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T)$? Justify your answer (you do not
 4037 have to explicitly compute the joint distribution). (1–2 sentences)
 4038 2 Assume that $p(\mathbf{x}_t|\mathbf{y}_1, \dots, \mathbf{y}_t) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$.
 4039 a) Compute $p(\mathbf{x}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_t)$
 4040 b) Compute $p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_t)$
 4041 c) At time $t+1$, we observe the value $\mathbf{y}_{t+1} = \hat{\mathbf{y}}$. Compute $p(\mathbf{x}_{t+1}|\mathbf{y}_1, \dots, \mathbf{y}_{t+1})$.

4042 6.4 Prove the relationship in Equation 6.44, which relates the standard definition
 4043 of the variance to the raw score expression for the variance.

4044 6.5 Prove the relationship in Equation 6.45, which relates the pairwise difference
 4045 between examples in a dataset with the raw score expression for the
 4046 variance.

4047 6.6 Express the Bernoulli distribution in the natural parameter form of the exponential
 4048 family (Equation (6.110)).

4049 6.7 Express the Binomial distribution as an exponential family distribution. Also
 4050 express the Beta distribution as an exponential family distribution. Show that
 4051 the product of the Beta and the Binomial distribution is also a member of
 4052 the exponential family.

4053 6.8 Derive the relationship in Section 6.5.2 in two ways:

- 4054 1 By completing the square
 4055 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

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

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

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

4056 Note that the normalizing constant c itself can be considered a (normalized)
 4057 Gaussian distribution either in \mathbf{a} or in \mathbf{b} with an “inflated” covariance matrix
 4058 $\mathbf{A} + \mathbf{B}$, 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:

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

4059 Here, $\mathbb{E}_x[x|y]$ denotes the expected value of x under the conditional distri-
 4060 bution $p(x|y)$.

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

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

4061 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 independ-
 4062 ent Gaussian noise. “Independent” implies that \mathbf{x} and \mathbf{w} are independent
 4063 random variables and that \mathbf{Q} is diagonal.

- 4064 1 Write down the likelihood $p(\mathbf{y}|\mathbf{x})$.
- 4065 2 The distribution $p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x}$ is Gaussian.¹ Compute the mean
 4066 $\boldsymbol{\mu}_y$ and the covariance $\boldsymbol{\Sigma}_y$. Derive your result in detail.
- 3 The random variable \mathbf{y} is being transformed according to the measure-
 ment mapping

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

4067 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
 4068 (measurement) noise.

- 4069 • Write down $p(\mathbf{z}|\mathbf{y})$.
 - 4070 • Compute $p(\mathbf{z})$, i.e., the mean $\boldsymbol{\mu}_z$ and the covariance $\boldsymbol{\Sigma}_z$. Derive your
 4071 result in detail.
- 4072 4 Now, a value $\hat{\mathbf{y}}$ is measured. Compute the posterior distribution $p(\mathbf{x}|\hat{\mathbf{y}})$.²
 4073 *Hint for solution:* Start by explicitly computing the joint Gaussian $p(\mathbf{x}, \mathbf{y})$.
 4074 This also requires to compute the cross-covariances $\text{Cov}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}, \mathbf{y}]$ and $\text{Cov}_{\mathbf{y}, \mathbf{x}}[\mathbf{y}, \mathbf{x}]$.
 4075 Then, apply the rules for Gaussian conditioning.

6.11 Probability integral transformation

4076 Given a continuous random variable x , with cdf $F_x(x)$. Show that the random
 4077 variable $y = F_x(x)$ is uniformly distributed.
 4078

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.