



In this chapter we discuss some specific examples of probability distributions and their properties. As well as being of interest in their own right, these distributions can form building blocks for more complex models and will be used extensively throughout the book.

One role for the distributions discussed in this chapter is to model the probability distribution  $p(\mathbf{x})$  of a random variable  $\mathbf{x}$ , given a finite set  $\mathbf{x}_1, \dots, \mathbf{x}_N$  of observations. This problem is known as *density estimation*. It should be emphasized that the problem of density estimation is fundamentally ill-posed, because there are infinitely many probability distributions that could have given rise to the observed finite data set. Indeed, any distribution  $p(\mathbf{x})$  that is non-zero at each of the data points  $\mathbf{x}_1, \dots, \mathbf{x}_N$  is a potential candidate. The issue of choosing an appropriate distribution relates to the problem of model selection, which has already been encountered in the context of polynomial curve fitting and which is a central issue in machine

### Section 1.2

learning.

We begin by considering distributions for discrete variables before exploring the Gaussian distribution for continuous variables. These are specific examples of *parametric* distributions, so called because they are governed by a relatively small number of adjustable parameters, such as the mean and variance of a Gaussian. To apply such models to the problem of density estimation, we need a procedure for determining suitable values for the parameters, given an observed data set, and our main focus will be on maximizing the likelihood function. In this chapter, we will assume that the data observations are independent and identically distributed (i.i.d.), whereas in future chapters we will explore more complex scenarios involving *structured data* where this assumption no longer holds.

One limitation of the parametric approach is that it assumes a specific functional form for the distribution, which may turn out to be inappropriate for a particular application. An alternative approach is given by *nonparametric* density estimation methods in which the form of the distribution typically depends on the size of the data set. Such models still contain parameters, but these control the model complexity rather than the form of the distribution. We end this chapter by briefly considering three nonparametric methods based respectively on histograms, nearest neighbours, and kernels. A major limitation of nonparametric techniques such as these is that they involve storing all the training data. In other words, the number of parameters grows with the size of the data set, so that the method become very inefficient for large data sets. Deep learning combines the efficiency of parametric models with the generality of nonparametric methods by considering flexible distributions based on neural networks having a large, but fixed, number of parameters.

### 3.1. Discrete Variables

---

We begin by considering simple distributions for discrete variables, starting with binary variables and then moving on to multi-state variables.

#### 3.1.1 Bernoulli distribution

Consider a single binary random variable  $x \in \{0, 1\}$ . For example,  $x$  might describe the outcome of flipping a coin, with  $x = 1$  representing ‘heads’ and  $x = 0$  representing ‘tails’. If this were a damaged coin, such as the one shown in [Figure 2.2](#), the probability of landing heads is not necessarily the same as that of landing tails. The probability of  $x = 1$  will be denoted by the parameter  $\mu$  so that

$$p(x = 1|\mu) = \mu \quad (3.1)$$

where  $0 \leq \mu \leq 1$ , from which it follows that  $p(x = 0|\mu) = 1 - \mu$ . The probability distribution over  $x$  can therefore be written in the form

$$\text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}, \quad (3.2)$$

#### Exercise 3.1

which is known as the *Bernoulli* distribution. It is easily verified that this distribution

is normalized and that it has mean and variance given by

$$\mathbb{E}[x] = \mu \quad (3.3)$$

$$\text{var}[x] = \mu(1 - \mu). \quad (3.4)$$

Now suppose we have a data set  $\mathcal{D} = \{x_1, \dots, x_N\}$  of observed values of  $x$ . We can construct the likelihood function, which is a function of  $\mu$ , on the assumption that the observations are drawn independently from  $p(x|\mu)$ , so that

$$p(\mathcal{D}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n}. \quad (3.5)$$

We can estimate a value for  $\mu$  by maximizing the likelihood function or equivalently by maximizing the logarithm of the likelihood, since the log is a monotonic function. The log likelihood function of the Bernoulli distribution is given by

$$\ln p(\mathcal{D}|\mu) = \sum_{n=1}^N \ln p(x_n|\mu) = \sum_{n=1}^N \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\}. \quad (3.6)$$

At this point, note that the log likelihood function depends on the  $N$  observations  $x_n$  only through their sum  $\sum_n x_n$ . This sum provides an example of a *sufficient statistic* for the data under this distribution. If we set the derivative of  $\ln p(\mathcal{D}|\mu)$  with respect to  $\mu$  equal to zero, we obtain the maximum likelihood estimator:

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n, \quad (3.7)$$

which is also known as the *sample mean*. Denoting the number of observations of  $x = 1$  (heads) within this data set by  $m$ , we can write (3.7) in the form

$$\mu_{\text{ML}} = \frac{m}{N} \quad (3.8)$$

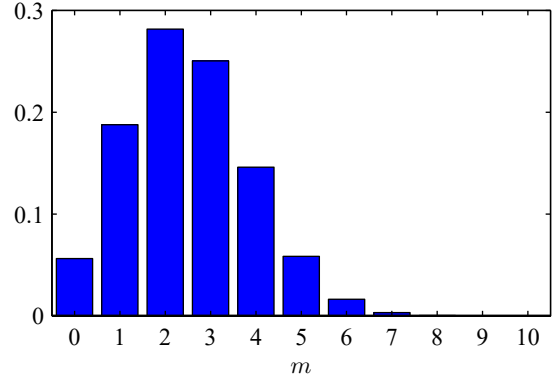
so that the probability of landing heads is given, in this maximum likelihood framework, by the fraction of observations of heads in the data set.

### 3.1.2 Binomial distribution

We can also work out the distribution for the binary variable  $x$  of the number  $m$  of observations of  $x = 1$ , given that the data set has size  $N$ . This is called the *binomial* distribution, and from (3.5) we see that it is proportional to  $\mu^m (1 - \mu)^{N-m}$ . To obtain the normalization coefficient, note that out of  $N$  coin flips, we have to add up all of the possible ways of obtaining  $m$  heads, so that the binomial distribution can be written as

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m} \quad (3.9)$$

**Figure 3.1** Histogram plot of the binomial distribution (3.9) as a function of  $m$  for  $N = 10$  and  $\mu = 0.25$ .



where

$$\binom{N}{m} \equiv \frac{N!}{(N-m)!m!} \quad (3.10)$$

is the number of ways of choosing  $m$  objects out of a total of  $N$  identical objects without replacement. Figure 3.1 shows a plot of the binomial distribution for  $N = 10$  and  $\mu = 0.25$ .

*Exercise 3.3*

The mean and variance of the binomial distribution can be found by using the results that, for independent events, the mean of the sum is the sum of the means and the variance of the sum is the sum of the variances. Because  $m = x_1 + \dots + x_N$  and because for each observation the mean and variance are given by (3.3) and (3.4), respectively, we have

*Exercise 2.10*

$$\mathbb{E}[m] \equiv \sum_{m=0}^N m \text{Bin}(m|N, \mu) = N\mu \quad (3.11)$$

$$\text{var}[m] \equiv \sum_{m=0}^N (m - \mathbb{E}[m])^2 \text{Bin}(m|N, \mu) = N\mu(1 - \mu). \quad (3.12)$$

*Exercise 3.4*

These results can also be proved directly by using calculus.

### 3.1.3 Multinomial distribution

Binary variables can be used to describe quantities that can take one of two possible values. Often, however, we encounter discrete variables that can take on one of  $K$  possible mutually exclusive states. Although there are various alternative ways to express such variables, we will see shortly that a particularly convenient representation is the 1-of- $K$  scheme, sometimes called ‘one-hot encoding’, in which the variable is represented by a  $K$ -dimensional vector  $\mathbf{x}$  in which one of the elements  $x_k$  equals 1 and all remaining elements equal 0. So, for instance, if we have a variable that can take  $K = 6$  states and a particular observation of the variable happens to

correspond to the state where  $x_3 = 1$ , then  $\mathbf{x}$  will be represented by

$$\mathbf{x} = (0, 0, 1, 0, 0)^T. \quad (3.13)$$

Note that such vectors satisfy  $\sum_{k=1}^K x_k = 1$ . If we denote the probability of  $x_k = 1$  by the parameter  $\mu_k$ , then the distribution of  $\mathbf{x}$  is given by

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad (3.14)$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T$ , and the parameters  $\mu_k$  are constrained to satisfy  $\mu_k \geq 0$  and  $\sum_k \mu_k = 1$ , because they represent probabilities. The distribution (3.14) can be regarded as a generalization of the Bernoulli distribution to more than two outcomes. It is easily seen that the distribution is normalized:

$$\sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) = \sum_{k=1}^K \mu_k = 1 \quad (3.15)$$

and that

$$\mathbb{E}[\mathbf{x}|\boldsymbol{\mu}] = \sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) \mathbf{x} = \boldsymbol{\mu}. \quad (3.16)$$

Now consider a data set  $\mathcal{D}$  of  $N$  independent observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . The corresponding likelihood function takes the form

$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k} \quad (3.17)$$

where we see that the likelihood function depends on the  $N$  data points only through the  $K$  quantities:

$$m_k = \sum_{n=1}^N x_{nk}, \quad (3.18)$$

which represent the number of observations of  $x_k = 1$ . These are called the *sufficient statistics* for this distribution. Note that the variables  $m_k$  are subject to the constraint

$$\sum_{k=1}^K m_k = N. \quad (3.19)$$

To find the maximum likelihood solution for  $\boldsymbol{\mu}$ , we need to maximize  $\ln p(\mathcal{D}|\boldsymbol{\mu})$  with respect to  $\mu_k$  taking account of the constraint (3.15) that the  $\mu_k$  must sum to one. This can be achieved using a Lagrange multiplier  $\lambda$  and maximizing

$$\sum_{k=1}^K m_k \ln \mu_k + \lambda \left( \sum_{k=1}^K \mu_k - 1 \right). \quad (3.20)$$

Section 3.4

Appendix C

Setting the derivative of (3.20) with respect to  $\mu_k$  to zero, we obtain

$$\mu_k = -m_k/\lambda. \quad (3.21)$$

We can solve for the Lagrange multiplier  $\lambda$  by substituting (3.21) into the constraint  $\sum_k \mu_k = 1$  to give  $\lambda = -N$ . Thus, we obtain the maximum likelihood solution for  $\mu_k$  in the form

$$\mu_k^{\text{ML}} = \frac{m_k}{N}, \quad (3.22)$$

which is the fraction of the  $N$  observations for which  $x_k = 1$ .

We can also consider the joint distribution of the quantities  $m_1, \dots, m_K$ , conditioned on the parameter vector  $\boldsymbol{\mu}$  and on the total number  $N$  of observations. From (3.17), this takes the form

$$\text{Mult}(m_1, m_2, \dots, m_K | \boldsymbol{\mu}, N) = \binom{N}{m_1 m_2 \dots m_K} \prod_{k=1}^K \mu_k^{m_k}, \quad (3.23)$$

which is known as the *multinomial* distribution. The normalization coefficient is the number of ways of partitioning  $N$  objects into  $K$  groups of size  $m_1, \dots, m_K$  and is given by

$$\binom{N}{m_1 m_2 \dots m_K} = \frac{N!}{m_1! m_2! \dots m_K!}. \quad (3.24)$$

Note that two-state quantities can be represented either as binary variables and modelled using the binomial distribution (3.9) or as 1-of-2 variables and modelled using the distribution (3.14) with  $K = 2$ .

## 3.2. The Multivariate Gaussian

### Section 2.3

The Gaussian, also known as the normal distribution, is a widely used model for the distribution of continuous variables. We have already seen that for of a single variable  $x$ , the Gaussian distribution can be written in the form

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} \quad (3.25)$$

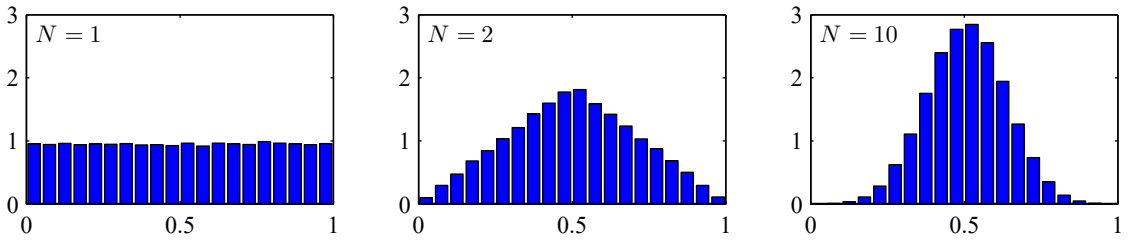
where  $\mu$  is the mean and  $\sigma^2$  is the variance. For a  $D$ -dimensional vector  $\mathbf{x}$ , the multivariate Gaussian distribution takes the form

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (3.26)$$

where  $\boldsymbol{\mu}$  is the  $D$ -dimensional mean vector,  $\boldsymbol{\Sigma}$  is the  $D \times D$  covariance matrix, and  $\det \boldsymbol{\Sigma}$  denotes the determinant of  $\boldsymbol{\Sigma}$ .

### Section 2.5

The Gaussian distribution arises in many different contexts and can be motivated from a variety of different perspectives. For example, we have already seen that for



**Figure 3.2** Histogram plots of the mean of  $N$  uniformly distributed numbers for various values of  $N$ . We observe that as  $N$  increases, the distribution tends towards a Gaussian.

### Exercise 3.8

a single real variable, the distribution that maximizes the entropy is the Gaussian. This property applies also to the multivariate Gaussian.

Another situation in which the Gaussian distribution arises is when we consider the sum of multiple random variables. The *central limit theorem* tells us that, subject to certain mild conditions, the sum of a set of random variables, which is of course itself a random variable, has a distribution that becomes increasingly Gaussian as the number of terms in the sum increases (Walker, 1969). We can illustrate this by considering  $N$  variables  $x_1, \dots, x_N$  each of which has a uniform distribution over the interval  $[0, 1]$  and then considering the distribution of the mean  $(x_1 + \dots + x_N)/N$ . For large  $N$ , this distribution tends to a Gaussian, as illustrated in Figure 3.2. In practice, the convergence to a Gaussian as  $N$  increases can be very rapid. One consequence of this result is that the binomial distribution (3.9), which is a distribution over  $m$  defined by the sum of  $N$  observations of the random binary variable  $x$ , will tend to a Gaussian as  $N \rightarrow \infty$  (see Figure 3.1 for  $N = 10$ ).

The Gaussian distribution has many important analytical properties, and we will consider several of these in detail. As a result, this section will be rather more technically involved than some of the earlier sections and will require familiarity with various matrix identities.

### Appendix A

#### 3.2.1 Geometry of the Gaussian

We begin by considering the geometrical form of the Gaussian distribution. The functional dependence of the Gaussian on  $\mathbf{x}$  is through the quadratic form

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}), \quad (3.27)$$

which appears in the exponent. The quantity  $\Delta$  is called the *Mahalanobis distance* from  $\boldsymbol{\mu}$  to  $\mathbf{x}$ . It reduces to the Euclidean distance when  $\boldsymbol{\Sigma}$  is the identity matrix. The Gaussian distribution is constant on surfaces in  $\mathbf{x}$ -space for which this quadratic form is constant.

### Exercise 3.11

First, note that the matrix  $\boldsymbol{\Sigma}$  can be taken to be symmetric, without loss of generality, because any antisymmetric component would disappear from the exponent. Now consider the eigenvector equation for the covariance matrix

$$\boldsymbol{\Sigma} \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad (3.28)$$

## Exercise 3.12

where  $i = 1, \dots, D$ . Because  $\Sigma$  is a real, symmetric matrix, its eigenvalues will be real, and its eigenvectors can be chosen to form an orthonormal set, so that

$$\mathbf{u}_i^T \mathbf{u}_j = I_{ij} \quad (3.29)$$

where  $I_{ij}$  is the  $i, j$  element of the identity matrix and satisfies

$$I_{ij} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise.} \end{cases} \quad (3.30)$$

## Exercise 3.13

The covariance matrix  $\Sigma$  can be expressed as an expansion in terms of its eigenvectors in the form

$$\Sigma = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T \quad (3.31)$$

and similarly the inverse covariance matrix  $\Sigma^{-1}$  can be expressed as

$$\Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T. \quad (3.32)$$

Substituting (3.32) into (3.27), the quadratic form becomes

$$\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad (3.33)$$

where we have defined

$$y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu}). \quad (3.34)$$

We can interpret  $\{y_i\}$  as a new coordinate system defined by the orthonormal vectors  $\mathbf{u}_i$  that are shifted and rotated with respect to the original  $x_i$  coordinates. Forming the vector  $\mathbf{y} = (y_1, \dots, y_D)^T$ , we have

$$\mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu}) \quad (3.35)$$

## Appendix A

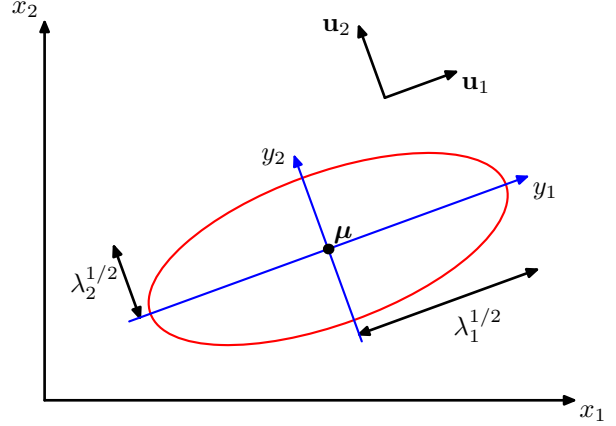
where  $\mathbf{U}$  is a matrix whose rows are given by  $\mathbf{u}_i^T$ . From (3.29) it follows that  $\mathbf{U}$  is an *orthogonal* matrix, i.e., it satisfies  $\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix.

The quadratic form, and hence the Gaussian density, is constant on surfaces for which (3.33) is constant. If all the eigenvalues  $\lambda_i$  are positive, then these surfaces represent ellipsoids, with their centres at  $\boldsymbol{\mu}$  and their axes oriented along  $\mathbf{u}_i$ , and with scaling factors in the directions of the axes given by  $\lambda_i^{1/2}$ , as illustrated in [Figure 3.3](#).

For the Gaussian distribution to be well defined, it is necessary for all the eigenvalues  $\lambda_i$  of the covariance matrix to be strictly positive, otherwise the distribution cannot be properly normalized. A matrix whose eigenvalues are strictly positive is said to be *positive definite*. When we discuss latent variable models, we will encounter Gaussian distributions for which one or more of the eigenvalues are zero, in



**Figure 3.3** The red curve shows the elliptical surface of constant probability density for a Gaussian in a two-dimensional space  $\mathbf{x} = (x_1, x_2)$  on which the density is  $\exp(-1/2)$  of its value at  $\mathbf{x} = \boldsymbol{\mu}$ . The axes of the ellipse are defined by the eigenvectors  $\mathbf{u}_i$  of the covariance matrix, with corresponding eigenvalues  $\lambda_i$ .



which case the distribution is singular and is confined to a subspace of lower dimensionality. If all the eigenvalues are non-negative, then the covariance matrix is said to be *positive semidefinite*.

Now consider the form of the Gaussian distribution in the new coordinate system defined by the  $y_i$ . In going from the  $\mathbf{x}$  to the  $\mathbf{y}$  coordinate system, we have a Jacobian matrix  $\mathbf{J}$  with elements given by

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ji} \quad (3.36)$$

where  $U_{ji}$  are the elements of the matrix  $\mathbf{U}^T$ . Using the orthonormality property of the matrix  $\mathbf{U}$ , we see that the square of the determinant of the Jacobian matrix is

$$|\mathbf{J}|^2 = |\mathbf{U}^T|^2 = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{I}| = 1 \quad (3.37)$$

and, hence,  $|\mathbf{J}| = 1$ . Also, the determinant  $|\boldsymbol{\Sigma}|$  of the covariance matrix can be written as the product of its eigenvalues, and hence

$$|\boldsymbol{\Sigma}|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2}. \quad (3.38)$$

Thus, in the  $y_j$  coordinate system, the Gaussian distribution takes the form

$$p(\mathbf{y}) = p(\mathbf{x}) |\mathbf{J}| = \prod_{j=1}^D \frac{1}{(2\pi\lambda_j)^{1/2}} \exp \left\{ -\frac{y_j^2}{2\lambda_j} \right\}, \quad (3.39)$$

which is the product of  $D$  independent univariate Gaussian distributions. The eigenvectors therefore define a new set of shifted and rotated coordinates with respect to which the joint probability distribution factorizes into a product of independent distributions. The integral of the distribution in the  $\mathbf{y}$  coordinate system is then

$$\int p(\mathbf{y}) \, d\mathbf{y} = \prod_{j=1}^D \int_{-\infty}^{\infty} \frac{1}{(2\pi\lambda_j)^{1/2}} \exp \left\{ -\frac{y_j^2}{2\lambda_j} \right\} \, dy_j = 1 \quad (3.40)$$

where we have used the result (2.51) for the normalization of the univariate Gaussian. This confirms that the multivariate Gaussian (3.26) is indeed normalized.

### 3.2.2 Moments

We now look at the moments of the Gaussian distribution and thereby provide an interpretation of the parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . The expectation of  $\mathbf{x}$  under the Gaussian distribution is given by

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \int \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\} \mathbf{x} d\mathbf{x} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \int \exp \left\{ -\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1} \mathbf{z} \right\} (\mathbf{z} + \boldsymbol{\mu}) d\mathbf{z}\end{aligned}\quad (3.41)$$

where we have changed variables using  $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$ . Note that the exponent is an even function of the components of  $\mathbf{z}$ , and because the integrals over these are taken over the range  $(-\infty, \infty)$ , the term in  $\mathbf{z}$  in the factor  $(\mathbf{z} + \boldsymbol{\mu})$  will vanish by symmetry. Thus,

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}, \quad (3.42)$$

and so we refer to  $\boldsymbol{\mu}$  as the mean of the Gaussian distribution.

We now consider second-order moments of the Gaussian. In the univariate case, we considered the second-order moment given by  $\mathbb{E}[x^2]$ . For the multivariate Gaussian, there are  $D^2$  second-order moments given by  $\mathbb{E}[x_i x_j]$ , which we can group together to form the matrix  $\mathbb{E}[\mathbf{x}\mathbf{x}^T]$ . This matrix can be written as

$$\begin{aligned}\mathbb{E}[\mathbf{x}\mathbf{x}^T] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \int \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\} \mathbf{x}\mathbf{x}^T d\mathbf{x} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \int \exp \left\{ -\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1} \mathbf{z} \right\} (\mathbf{z} + \boldsymbol{\mu})(\mathbf{z} + \boldsymbol{\mu})^T d\mathbf{z}\end{aligned}\quad (3.43)$$

where again we have changed variables using  $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$ . Note that the cross-terms involving  $\boldsymbol{\mu}\mathbf{z}^T$  and  $\boldsymbol{\mu}^T \mathbf{z}$  will again vanish by symmetry. The term  $\boldsymbol{\mu}\boldsymbol{\mu}^T$  is constant and can be taken outside the integral, which itself is unity because the Gaussian distribution is normalized. Consider the term involving  $\mathbf{z}\mathbf{z}^T$ . Again, we can make use of the eigenvector expansion of the covariance matrix given by (3.28), together with the completeness of the set of eigenvectors, to write

$$\mathbf{z} = \sum_{j=1}^D y_j \mathbf{u}_j \quad (3.44)$$

where  $y_j = \mathbf{u}_j^T \mathbf{z}$ , which gives

$$\begin{aligned}
 & \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp \left\{ -\frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z} \right\} \mathbf{z} \mathbf{z}^T d\mathbf{z} \\
 &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \sum_{i=1}^D \sum_{j=1}^D \mathbf{u}_i \mathbf{u}_j^T \int \exp \left\{ -\sum_{k=1}^D \frac{y_k^2}{2\lambda_k} \right\} y_i y_j d\mathbf{y} \\
 &= \sum_{i=1}^D \mathbf{u}_i \mathbf{u}_i^T \lambda_i = \Sigma
 \end{aligned} \tag{3.45}$$

where we have made use of the eigenvector equation (3.28), together with the fact that the integral on the middle line vanishes by symmetry unless  $i = j$ . In the final line we have made use of the results (2.53) and (3.38), together with (3.31). Thus, we have

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \Sigma. \tag{3.46}$$

When defining the variance for a single random variable, we subtracted the mean before taking the second moment. Similarly, in the multivariate case it is again convenient to subtract off the mean, giving rise to the *covariance* of a random vector  $\mathbf{x}$  defined by

$$\text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T]. \tag{3.47}$$

For the specific case of a Gaussian distribution, we can make use of  $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$ , together with the result (3.46), to give

$$\text{cov}[\mathbf{x}] = \Sigma. \tag{3.48}$$

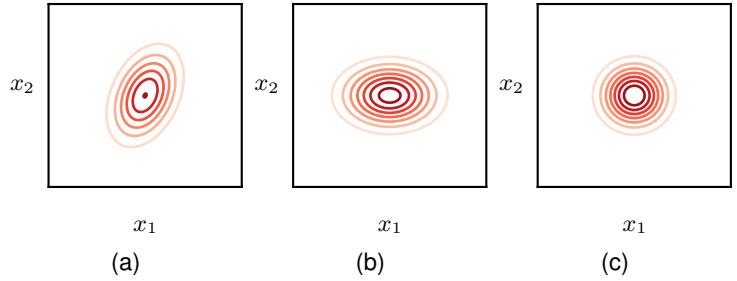
Because the parameter matrix  $\Sigma$  governs the covariance of  $\mathbf{x}$  under the Gaussian distribution, it is called the covariance matrix.

### 3.2.3 Limitations

Although the Gaussian distribution (3.26) is often used as a simple density model, it suffers from some significant limitations. Consider the number of free parameters in the distribution. A general symmetric covariance matrix  $\Sigma$  will have  $D(D+1)/2$  independent parameters, and there are another  $D$  independent parameters in  $\boldsymbol{\mu}$ , giving  $D(D+3)/2$  parameters in total. For large  $D$ , the total number of parameters therefore grows quadratically with  $D$ , and the computational task of manipulating and inverting the large matrices can become prohibitive. One way to address this problem is to use restricted forms of the covariance matrix. If we consider covariance matrices that are *diagonal*, so that  $\Sigma = \text{diag}(\sigma_i^2)$ , we then have a total of  $2D$  independent parameters in the density model. The corresponding contours of constant density are given by axis-aligned ellipsoids. We could further restrict the covariance matrix to be proportional to the identity matrix,  $\Sigma = \sigma^2 \mathbf{I}$ , known as an *isotropic* covariance, giving  $D+1$  independent parameters in the model together with spherical surfaces of constant density. The three possibilities of general, diagonal, and isotropic covariance matrices are illustrated in Figure 3.4. Unfortunately,

*Exercise 3.15*

**Figure 3.4** Contours of constant probability density for a Gaussian distribution in two dimensions in which the covariance matrix is (a) of general form, (b) diagonal, in which case the elliptical contours are aligned with the coordinate axes, and (c) proportional to the identity matrix, in which case the contours are concentric circles.



whereas such approaches limit the number of degrees of freedom in the distribution and make inversion of the covariance matrix a much faster operation, they also greatly restrict the form of the probability density and limit its ability to capture interesting correlations in the data.

A further limitation of the Gaussian distribution is that it is intrinsically unimodal (i.e., has a single maximum) and so is unable to provide a good approximation to multimodal distributions. Thus, the Gaussian distribution can be both too flexible, in the sense of having too many parameters, and too limited in the range of distributions that it can adequately represent. We will see later that the introduction of *latent* variables, also called *hidden* variables or *unobserved* variables, allows both of these problems to be addressed. In particular, a rich family of multimodal distributions is obtained by introducing discrete latent variables leading to mixtures of Gaussians. Similarly, the introduction of continuous latent variables leads to models in which the number of free parameters can be controlled independently of the dimensionality  $D$  of the data space while still allowing the model to capture the dominant correlations in the data set.

Section 3.2.9

Chapter 16

### 3.2.4 Conditional distribution

An important property of a multivariate Gaussian distribution is that if two sets of variables are jointly Gaussian, then the conditional distribution of one set conditioned on the other is again Gaussian. Similarly, the marginal distribution of either set is also Gaussian.

First, consider the case of conditional distributions. Suppose that  $\mathbf{x}$  is a  $D$ -dimensional vector with Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and that we partition  $\mathbf{x}$  into two disjoint subsets  $\mathbf{x}_a$  and  $\mathbf{x}_b$ . Without loss of generality, we can take  $\mathbf{x}_a$  to form the first  $M$  components of  $\mathbf{x}$ , with  $\mathbf{x}_b$  comprising the remaining  $D - M$  components, so that

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}. \quad (3.49)$$

We also define corresponding partitions of the mean vector  $\boldsymbol{\mu}$  given by

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad (3.50)$$

and of the covariance matrix  $\Sigma$  given by

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}. \quad (3.51)$$

Note that the symmetry  $\Sigma^T = \Sigma$  of the covariance matrix implies that  $\Sigma_{aa}$  and  $\Sigma_{bb}$  are symmetric and that  $\Sigma_{ba} = \Sigma_{ab}^T$ .

In many situations, it will be convenient to work with the inverse of the covariance matrix:

$$\Lambda \equiv \Sigma^{-1}, \quad (3.52)$$

which is known as the *precision matrix*. In fact, we will see that some properties of Gaussian distributions are most naturally expressed in terms of the covariance, whereas others take a simpler form when viewed in terms of the precision. We therefore also introduce the partitioned form of the precision matrix:

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} \quad (3.53)$$

corresponding to the partitioning (3.49) of the vector  $\mathbf{x}$ . Because the inverse of a symmetric matrix is also symmetric, we see that  $\Lambda_{aa}$  and  $\Lambda_{bb}$  are symmetric and that  $\Lambda_{ba} = \Lambda_{ab}^T$ . It should be stressed at this point that, for instance,  $\Lambda_{aa}$  is not simply given by the inverse of  $\Sigma_{aa}$ . In fact, we will shortly examine the relation between the inverse of a partitioned matrix and the inverses of its partitions.

### Exercise 3.16

We begin by finding an expression for the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$ . From the product rule of probability, we see that this conditional distribution can be evaluated from the joint distribution  $p(\mathbf{x}) = p(\mathbf{x}_a, \mathbf{x}_b)$  simply by fixing  $\mathbf{x}_b$  to the observed value and normalizing the resulting expression to obtain a valid probability distribution over  $\mathbf{x}_a$ . Instead of performing this normalization explicitly, we can obtain the solution more efficiently by considering the quadratic form in the exponent of the Gaussian distribution given by (3.27) and then reinstating the normalization coefficient at the end of the calculation. If we make use of the partitioning (3.49), (3.50), and (3.53), we obtain

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) = & \\ & -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \Lambda_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \Lambda_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ & -\frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \Lambda_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \Lambda_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned} \quad (3.54)$$

We see that as a function of  $\mathbf{x}_a$ , this is again a quadratic form, and hence, the corresponding conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  will be Gaussian. Because this distribution is completely characterized by its mean and its covariance, our goal will be to identify expressions for the mean and covariance of  $p(\mathbf{x}_a|\mathbf{x}_b)$  by inspection of (3.54).

This is an example of a rather common operation associated with Gaussian distributions, sometimes called ‘completing the square’, in which we are given a

quadratic form defining the exponent terms in a Gaussian distribution and we need to determine the corresponding mean and covariance. Such problems can be solved straightforwardly by noting that the exponent in a general Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  can be written as

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const} \quad (3.55)$$

where ‘const’ denotes terms that are independent of  $\mathbf{x}$ . We have also made use of the symmetry of  $\boldsymbol{\Sigma}$ . Thus, if we take our general quadratic form and express it in the form given by the right-hand side of (3.55), then we can immediately equate the matrix of coefficients entering the second-order term in  $\mathbf{x}$  to the inverse covariance matrix  $\boldsymbol{\Sigma}^{-1}$  and the coefficient of the linear term in  $\mathbf{x}$  to  $\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}$ , from which we can obtain  $\boldsymbol{\mu}$ .

Now let us apply this procedure to the conditional Gaussian distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  for which the quadratic form in the exponent is given by (3.54). We will denote the mean and covariance of this distribution by  $\boldsymbol{\mu}_{a|b}$  and  $\boldsymbol{\Sigma}_{a|b}$ , respectively. Consider the functional dependence of (3.54) on  $\mathbf{x}_a$  in which  $\mathbf{x}_b$  is regarded as a constant. If we pick out all terms that are second order in  $\mathbf{x}_a$ , we have

$$-\frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{aa}\mathbf{x}_a \quad (3.56)$$

from which we can immediately conclude that the covariance (inverse precision) of  $p(\mathbf{x}_a|\mathbf{x}_b)$  is given by

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1}. \quad (3.57)$$

Now consider all the terms in (3.54) that are linear in  $\mathbf{x}_a$ :

$$\mathbf{x}_a^T \{ \boldsymbol{\Lambda}_{aa}\boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \} \quad (3.58)$$

where we have used  $\boldsymbol{\Lambda}_{ba}^T = \boldsymbol{\Lambda}_{ab}$ . From our discussion of the general form (3.55), the coefficient of  $\mathbf{x}_a$  in this expression must equal  $\boldsymbol{\Sigma}_{a|b}^{-1}\boldsymbol{\mu}_{a|b}$  and, hence,

$$\begin{aligned} \boldsymbol{\mu}_{a|b} &= \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa}\boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned} \quad (3.59)$$

where we have made use of (3.57).

The results (3.57) and (3.59) are expressed in terms of the partitioned precision matrix of the original joint distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$ . We can also express these results in terms of the corresponding partitioned covariance matrix. To do this, we make use of the following identity for the inverse of a partitioned matrix:

*Exercise 3.18*

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix} \quad (3.60)$$

where we have defined

$$\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}. \quad (3.61)$$

The quantity  $\mathbf{M}^{-1}$  is known as the *Schur complement* of the matrix on the left-hand side of (3.60) with respect to the submatrix  $\mathbf{D}$ . Using the definition

$$\begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} \quad (3.62)$$

and making use of (3.60), we have

$$\Lambda_{aa} = (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1} \quad (3.63)$$

$$\Lambda_{ab} = -(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}. \quad (3.64)$$

From these we obtain the following expressions for the mean and covariance of the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$ :

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \Sigma_{ab}\Sigma_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (3.65)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}. \quad (3.66)$$

Comparing (3.57) and (3.66), we see that the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  takes a simpler form when expressed in terms of the partitioned precision matrix than when it is expressed in terms of the partitioned covariance matrix. Note that the mean of the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$ , given by (3.65), is a linear function of  $\mathbf{x}_b$  and that the covariance, given by (3.66), is independent of  $\mathbf{x}_b$ . This represents an example of a *linear-Gaussian* model.

Section 11.1.4

### 3.2.5 Marginal distribution

We have seen that if a joint distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$  is Gaussian, then the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  will again be Gaussian. Now we turn to a discussion of the marginal distribution given by

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b, \quad (3.67)$$

which, as we will see, is also Gaussian. Once again, our strategy for calculating this distribution will be to focus on the quadratic form in the exponent of the joint distribution and thereby to identify the mean and covariance of the marginal distribution  $p(\mathbf{x}_a)$ .

The quadratic form for the joint distribution can be expressed, using the partitioned precision matrix, in the form (3.54). Our goal is to integrate out  $\mathbf{x}_b$ , which is most easily achieved by first considering the terms involving  $\mathbf{x}_b$  and then completing the square to facilitate the integration. Picking out just those terms that involve  $\mathbf{x}_b$ , we have

$$-\frac{1}{2}\mathbf{x}_b^T \Lambda_{bb} \mathbf{x}_b + \mathbf{x}_b^T \mathbf{m} = -\frac{1}{2}(\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m}) + \frac{1}{2}\mathbf{m}^T \Lambda_{bb}^{-1} \mathbf{m} \quad (3.68)$$

where we have defined

$$\mathbf{m} = \Lambda_{bb}\boldsymbol{\mu}_b - \Lambda_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a). \quad (3.69)$$

We see that the dependence on  $\mathbf{x}_b$  has been cast into the standard quadratic form of a Gaussian distribution corresponding to the first term on the right-hand side of (3.68) plus a term that does not depend on  $\mathbf{x}_b$  (but that does depend on  $\mathbf{x}_a$ ). Thus, when we take the exponential of this quadratic form, we see that the integration over  $\mathbf{x}_b$  required by (3.67) will take the form

$$\int \exp \left\{ -\frac{1}{2} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m}) \right\} d\mathbf{x}_b. \quad (3.70)$$

This integration is easily performed by noting that it is the integral over an unnormalized Gaussian, and so the result will be the reciprocal of the normalization coefficient. We know from the form of the normalized Gaussian given by (3.26) that this coefficient is independent of the mean and depends only on the determinant of the covariance matrix. Thus, by completing the square with respect to  $\mathbf{x}_b$ , we can integrate out  $\mathbf{x}_b$  so that the only term remaining from the contributions on the left-hand side of (3.68) that depends on  $\mathbf{x}_a$  is the last term on the right-hand side of (3.68) in which  $\mathbf{m}$  is given by (3.69). Combining this term with the remaining terms from (3.54) that depend on  $\mathbf{x}_a$ , we obtain

$$\begin{aligned} & \frac{1}{2} [\Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a)]^T \Lambda_{bb}^{-1} [\Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a)] \\ & - \frac{1}{2} \mathbf{x}_a^T \Lambda_{aa} \mathbf{x}_a + \mathbf{x}_a^T (\Lambda_{aa} \boldsymbol{\mu}_a + \Lambda_{ab} \boldsymbol{\mu}_b) + \text{const} \\ & = -\frac{1}{2} \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mathbf{x}_a \\ & \quad + \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \boldsymbol{\mu}_a + \text{const} \end{aligned} \quad (3.71)$$

where ‘const’ denotes quantities independent of  $\mathbf{x}_a$ . Again, by comparison with (3.55), we see that the covariance of the marginal distribution  $p(\mathbf{x}_a)$  is given by

$$\Sigma_a = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1}. \quad (3.72)$$

Similarly, the mean is given by

$$\Sigma_a (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \boldsymbol{\mu}_a = \boldsymbol{\mu}_a \quad (3.73)$$

where we have used (3.72). The covariance (3.72) is expressed in terms of the partitioned precision matrix given by (3.53). We can rewrite this in terms of the corresponding partitioning of the covariance matrix given by (3.51), as we did for the conditional distribution. These partitioned matrices are related by

$$\begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}. \quad (3.74)$$

Making use of (3.60), we then have

$$(\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1} = \Sigma_{aa}. \quad (3.75)$$



Thus, we obtain the intuitively satisfying result that the marginal distribution  $p(\mathbf{x}_a)$  has mean and covariance given by

$$\mathbb{E}[\mathbf{x}_a] = \boldsymbol{\mu}_a \quad (3.76)$$

$$\text{cov}[\mathbf{x}_a] = \boldsymbol{\Sigma}_{aa}. \quad (3.77)$$

We see that for a marginal distribution, the mean and covariance are most simply expressed in terms of the partitioned covariance matrix, in contrast to the conditional distribution for which the partitioned precision matrix gives rise to simpler expressions.

Our results for the marginal and conditional distributions of a partitioned Gaussian can be summarized as follows. Given a joint Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with  $\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$  and the following partitions

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad (3.78)$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix} \quad (3.79)$$

then the conditional distribution is given by

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \quad (3.80)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (3.81)$$

and the marginal distribution is given by

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}). \quad (3.82)$$

We illustrate the idea of conditional and marginal distributions associated with a multivariate Gaussian using an example involving two variables in [Figure 3.5](#).

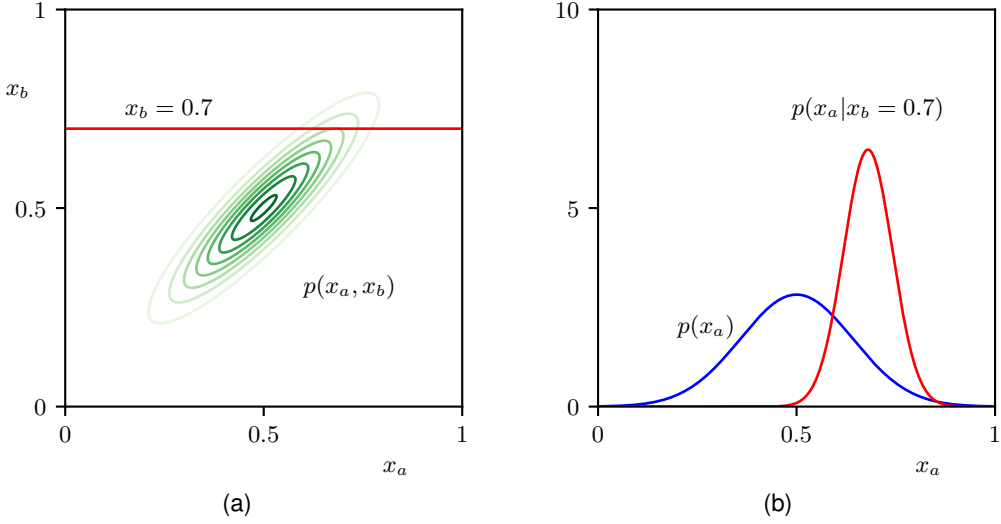
### 3.2.6 Bayes' theorem

In Sections 3.2.4 and 3.2.5 we considered a Gaussian  $p(\mathbf{x})$  in which we partitioned the vector  $\mathbf{x}$  into two subvectors  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$  and then found expressions for the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  and the marginal distribution  $p(\mathbf{x}_a)$ . We noted that the mean of the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  was a linear function of  $\mathbf{x}_b$ . Here we will suppose that we are given a Gaussian marginal distribution  $p(\mathbf{x})$  and a Gaussian conditional distribution  $p(\mathbf{y}|\mathbf{x})$  in which  $p(\mathbf{y}|\mathbf{x})$  has a mean that is a linear function of  $\mathbf{x}$  and a covariance that is independent of  $\mathbf{x}$ . This is an example of a *linear-Gaussian model* (Roweis and Ghahramani, 1999). We wish to find the marginal distribution  $p(\mathbf{y})$  and the conditional distribution  $p(\mathbf{x}|\mathbf{y})$ . This is a structure that arises in several types of generative model and it will prove convenient to derive the general results here.

We will take the marginal and conditional distributions to be

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (3.83)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \quad (3.84)$$



**Figure 3.5** (a) Contours of a Gaussian distribution  $p(x_a, x_b)$  over two variables. (b) The marginal distribution  $p(x_a)$  (blue curve) and the conditional distribution  $p(x_a|x_b)$  for  $x_b = 0.7$  (red curve).

where  $\boldsymbol{\mu}$ ,  $\mathbf{A}$ , and  $\mathbf{b}$  are parameters governing the means, and  $\boldsymbol{\Lambda}$  and  $\mathbf{L}$  are precision matrices. If  $\mathbf{x}$  has dimensionality  $M$  and  $\mathbf{y}$  has dimensionality  $D$ , then the matrix  $\mathbf{A}$  has size  $D \times M$ .

First we find an expression for the joint distribution over  $\mathbf{x}$  and  $\mathbf{y}$ . To do this, we define

$$\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \quad (3.85)$$

and then consider the log of the joint distribution:

$$\begin{aligned} \ln p(\mathbf{z}) &= \ln p(\mathbf{x}) + \ln p(\mathbf{y}|\mathbf{x}) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}) \\ &\quad -\frac{1}{2}(\mathbf{y} - \mathbf{Ax} - \mathbf{b})^T \mathbf{L} (\mathbf{y} - \mathbf{Ax} - \mathbf{b}) + \text{const} \end{aligned} \quad (3.86)$$

where ‘const’ denotes terms independent of  $\mathbf{x}$  and  $\mathbf{y}$ . As before, we see that this is a quadratic function of the components of  $\mathbf{z}$ , and hence,  $p(\mathbf{z})$  is Gaussian distribution. To find the precision of this Gaussian, we consider the second-order terms in (3.86), which can be written as

$$\begin{aligned} &-\frac{1}{2}\mathbf{x}^T (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}) \mathbf{x} - \frac{1}{2}\mathbf{y}^T \mathbf{L} \mathbf{y} + \frac{1}{2}\mathbf{y}^T \mathbf{L} \mathbf{A} \mathbf{x} + \frac{1}{2}\mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{y} \\ &= -\frac{1}{2} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = -\frac{1}{2} \mathbf{z}^T \mathbf{R} \mathbf{z} \end{aligned} \quad (3.87)$$

and so the Gaussian distribution over  $\mathbf{z}$  has precision (inverse covariance) matrix

given by

$$\mathbf{R} = \begin{pmatrix} \mathbf{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{pmatrix}. \quad (3.88)$$

The covariance matrix is found by taking the inverse of the precision, which can be done using the matrix inversion formula (3.60) to give

*Exercise 3.23*

$$\text{cov}[\mathbf{z}] = \mathbf{R}^{-1} = \begin{pmatrix} \mathbf{\Lambda}^{-1} & \mathbf{\Lambda}^{-1} \mathbf{A}^T \\ \mathbf{A} \mathbf{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^T \end{pmatrix}. \quad (3.89)$$

Similarly, we can find the mean of the Gaussian distribution over  $\mathbf{z}$  by identifying the linear terms in (3.86), which are given by

$$\mathbf{x}^T \mathbf{\Lambda} \boldsymbol{\mu} - \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{b} + \mathbf{y}^T \mathbf{L} \mathbf{b} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \mathbf{\Lambda} \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix}. \quad (3.90)$$

Using our earlier result (3.55) obtained by completing the square over the quadratic form of a multivariate Gaussian, we find that the mean of  $\mathbf{z}$  is given by

$$\mathbb{E}[\mathbf{z}] = \mathbf{R}^{-1} \begin{pmatrix} \mathbf{\Lambda} \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix}. \quad (3.91)$$

*Exercise 3.24*

Making use of (3.89), we then obtain

$$\mathbb{E}[\mathbf{z}] = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \end{pmatrix}. \quad (3.92)$$

Next we find an expression for the marginal distribution  $p(\mathbf{y})$  in which we have marginalized over  $\mathbf{x}$ . Recall that the marginal distribution over a subset of the components of a Gaussian random vector takes a particularly simple form when expressed in terms of the partitioned covariance matrix. Specifically, its mean and covariance are given by (3.76) and (3.77), respectively. Making use of (3.89) and (3.92), we see that the mean and covariance of the marginal distribution  $p(\mathbf{y})$  are given by

*Section 3.2*

$$\mathbb{E}[\mathbf{y}] = \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \quad (3.93)$$

$$\text{cov}[\mathbf{y}] = \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^T. \quad (3.94)$$

A special case of this result is when  $\mathbf{A} = \mathbf{I}$ , in which case the marginal distribution reduces to the convolution of two Gaussians, for which we see that the mean of the convolution is the sum of the means of the two Gaussians and the covariance of the convolution is the sum of their covariances.

Finally, we seek an expression for the conditional  $p(\mathbf{x}|\mathbf{y})$ . Recall that the results for the conditional distribution are most easily expressed in terms of the partitioned precision matrix, using (3.57) and (3.59). Applying these results to (3.89) and (3.92), we see that the conditional distribution  $p(\mathbf{x}|\mathbf{y})$  has mean and covariance given by

*Section 3.2*

$$\mathbb{E}[\mathbf{x}|\mathbf{y}] = (\mathbf{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1} \{ \mathbf{A}^T \mathbf{L} (\mathbf{y} - \mathbf{b}) + \mathbf{\Lambda} \boldsymbol{\mu} \} \quad (3.95)$$

$$\text{cov}[\mathbf{x}|\mathbf{y}] = (\mathbf{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}. \quad (3.96)$$

The evaluation of this conditional distribution can be seen as an example of Bayes' theorem, in which we interpret  $p(\mathbf{x})$  as a prior distribution over  $\mathbf{x}$ . If the variable  $\mathbf{y}$  is observed, then the conditional distribution  $p(\mathbf{x}|\mathbf{y})$  represents the corresponding posterior distribution over  $\mathbf{x}$ . Having found the marginal and conditional distributions, we have effectively expressed the joint distribution  $p(\mathbf{z}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$  in the form  $p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$ .

These results can be summarized as follows. Given a marginal Gaussian distribution for  $\mathbf{x}$  and a conditional Gaussian distribution for  $\mathbf{y}$  given  $\mathbf{x}$  in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (3.97)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}), \quad (3.98)$$

then the marginal distribution of  $\mathbf{y}$  and the conditional distribution of  $\mathbf{x}$  given  $\mathbf{y}$  are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \quad (3.99)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad (3.100)$$

where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}. \quad (3.101)$$

### 3.2.7 Maximum likelihood

Given a data set  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  in which the observations  $\{\mathbf{x}_n\}$  are assumed to be drawn independently from a multivariate Gaussian distribution, we can estimate the parameters of the distribution by maximum likelihood. The log likelihood function is given by

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}). \quad (3.102)$$

By simple rearrangement, we see that the likelihood function depends on the data set only through the two quantities

$$\sum_{n=1}^N \mathbf{x}_n, \quad \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T. \quad (3.103)$$

These are known as the *sufficient statistics* for the Gaussian distribution. Using (A.19), the derivative of the log likelihood with respect to  $\boldsymbol{\mu}$  is given by

$$\frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}), \quad (3.104)$$

and setting this derivative to zero, we obtain the solution for the maximum likelihood estimate of the mean:

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad (3.105)$$

which is the mean of the observed set of data points. The maximization of (3.102) with respect to  $\Sigma$  is rather more involved. The simplest approach is to ignore the symmetry constraint and show that the resulting solution is symmetric as required. Alternative derivations of this result, which impose the symmetry and positive definiteness constraints explicitly, can be found in Magnus and Neudecker (1999). The result is as expected and takes the form

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T, \quad (3.106)$$

which involves  $\boldsymbol{\mu}_{\text{ML}}$  because this is the result of a joint maximization with respect to  $\boldsymbol{\mu}$  and  $\Sigma$ . Note that the solution (3.105) for  $\boldsymbol{\mu}_{\text{ML}}$  does not depend on  $\Sigma_{\text{ML}}$ , and so we can first evaluate  $\boldsymbol{\mu}_{\text{ML}}$  and then use this to evaluate  $\Sigma_{\text{ML}}$ .

If we evaluate the expectations of the maximum likelihood solutions under the true distribution, we obtain the following results

$$\mathbb{E}[\boldsymbol{\mu}_{\text{ML}}] = \boldsymbol{\mu} \quad (3.107)$$

$$\mathbb{E}[\Sigma_{\text{ML}}] = \frac{N-1}{N} \Sigma. \quad (3.108)$$

We see that the expectation of the maximum likelihood estimate for the mean is equal to the true mean. However, the maximum likelihood estimate for the covariance has an expectation that is less than the true value, and hence, it is biased. We can correct this bias by defining a different estimator  $\tilde{\Sigma}$  given by

$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T. \quad (3.109)$$

Clearly from (3.106) and (3.108), the expectation of  $\tilde{\Sigma}$  is equal to  $\Sigma$ .

### 3.2.8 Sequential estimation

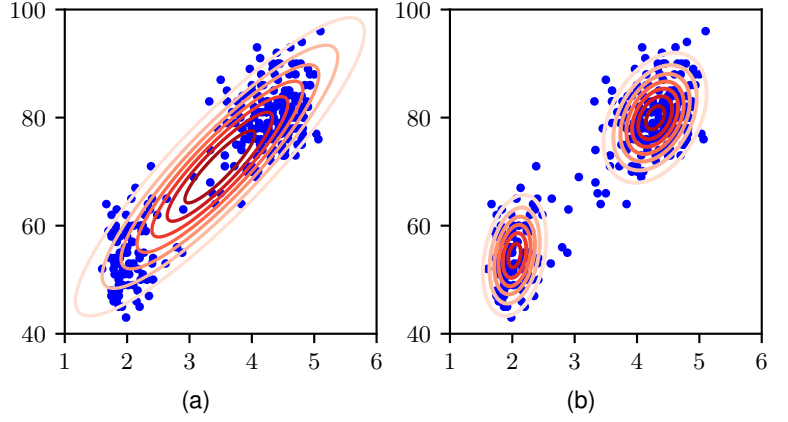
Our discussion of the maximum likelihood solution represents a *batch* method in which the entire training data set is considered at once. An alternative is to use *sequential* methods, which allow data points to be processed one at a time and then discarded. These are important for online applications and for large data when the batch processing of all data points at once is infeasible.

Consider the result (3.105) for the maximum likelihood estimator of the mean  $\boldsymbol{\mu}_{\text{ML}}$ , which we will denote by  $\boldsymbol{\mu}_{\text{ML}}^{(N)}$  when it is based on  $N$  observations. If we

*Exercise 3.28*

*Exercise 3.29*

**Figure 3.6** Plots of the Old Faithful data in which the red curves are contours of constant probability density. (a) A single Gaussian distribution which has been fitted to the data using maximum likelihood. Note that this distribution fails to capture the two clumps in the data and indeed places much of its probability mass in the central region between the clumps where the data are relatively sparse. (b) The distribution given by a linear combination of two Gaussians, also fitted by maximum likelihood, which gives a better representation of the data.



dissect out the contribution from the final data point  $\mathbf{x}_N$ , we obtain

$$\begin{aligned}
 \boldsymbol{\mu}_{\text{ML}}^{(N)} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\
 &= \frac{1}{N} \mathbf{x}_N + \frac{1}{N} \sum_{n=1}^{N-1} \mathbf{x}_n \\
 &= \frac{1}{N} \mathbf{x}_N + \frac{N-1}{N} \boldsymbol{\mu}_{\text{ML}}^{(N-1)} \\
 &= \boldsymbol{\mu}_{\text{ML}}^{(N-1)} + \frac{1}{N} (\mathbf{x}_N - \boldsymbol{\mu}_{\text{ML}}^{(N-1)}). \tag{3.110}
 \end{aligned}$$

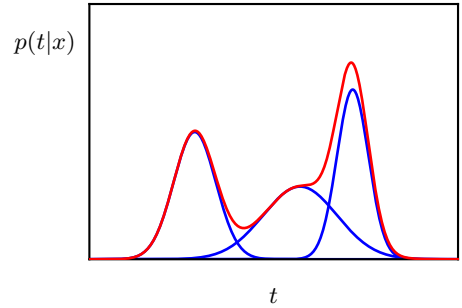
This result has a nice interpretation, as follows. After observing  $N-1$  data points, we estimate  $\boldsymbol{\mu}$  by  $\boldsymbol{\mu}_{\text{ML}}^{(N-1)}$ . We now observe data point  $\mathbf{x}_N$ , and we obtain our revised estimate  $\boldsymbol{\mu}_{\text{ML}}^{(N)}$  by moving the old estimate a small amount, proportional to  $1/N$ , in the direction of the ‘error signal’  $(\mathbf{x}_N - \boldsymbol{\mu}_{\text{ML}}^{(N-1)})$ . Note that, as  $N$  increases, so the contributions from successive data points get smaller.

### 3.2.9 Mixtures of Gaussians

Although the Gaussian distribution has some important analytical properties, it suffers from significant limitations when used to model real data sets. Consider the example shown in Figure 3.6(a). This is known as the ‘Old Faithful’ data set, and comprises 272 measurements of the eruption of the Old Faithful geyser in Yellowstone National Park in the USA. Each measurement gives the duration of the eruption in minutes (horizontal axis) and the time in minutes to the next eruption (vertical axis). We see that the data set forms two dominant clumps, and that a simple Gaussian distribution is unable to capture this structure.

We might expect that a superposition of two Gaussian distributions would be able to do a much better job of representing the structure in this data set, and indeed

**Figure 3.7** Example of a Gaussian mixture distribution in one dimension showing three Gaussians (each scaled by a coefficient) in blue and their sum in red.



this proves to be the case, as can be seen from Figure 3.6(b). Such superpositions, formed by taking linear combinations of more basic distributions such as Gaussians, can be formulated as probabilistic models known as *mixture distributions*. In this section we will consider Gaussians to illustrate the framework of mixture models. More generally, mixture models can comprise linear combinations of other distributions, for example mixtures of Bernoulli distributions for binary variables. In Figure 3.7 we see that a linear combination of Gaussians can give rise to very complex densities. By using a sufficient number of Gaussians and by adjusting their means and covariances as well as the coefficients in the linear combination, almost any continuous distribution can be approximated to arbitrary accuracy.

We therefore consider a superposition of  $K$  Gaussian densities of the form

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (3.111)$$

which is called a *mixture of Gaussians*. Each Gaussian density  $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  is called a *component* of the mixture and has its own mean  $\boldsymbol{\mu}_k$  and covariance  $\boldsymbol{\Sigma}_k$ . Contour and surface plots for a Gaussian mixture in two dimensions having three components are shown in Figure 3.8.

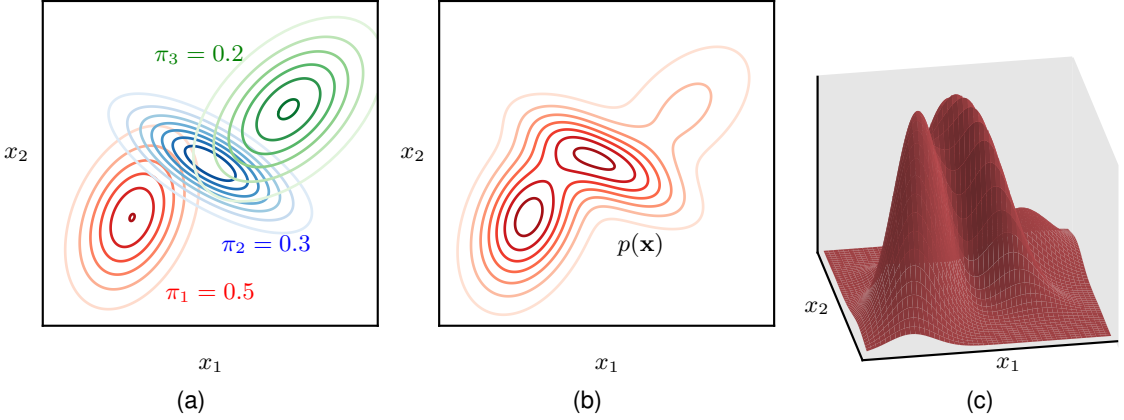
The parameters  $\pi_k$  in (3.111) are called *mixing coefficients*. If we integrate both sides of (3.111) with respect to  $\mathbf{x}$ , and note that both  $p(\mathbf{x})$  and the individual Gaussian components are normalized, we obtain

$$\sum_{k=1}^K \pi_k = 1. \quad (3.112)$$

Also, given that  $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \geq 0$ , a sufficient condition for the requirement  $p(\mathbf{x}) \geq 0$  is that  $\pi_k \geq 0$  for all  $k$ . Combining this with the condition (3.112), we obtain

$$0 \leq \pi_k \leq 1. \quad (3.113)$$

We can therefore see that the mixing coefficients satisfy the requirements to be probabilities, and we will show that this probabilistic interpretation of mixture distributions is very powerful.



**Figure 3.8** Illustration of a mixture of three Gaussians in a two-dimensional space. (a) Contours of constant density for each of the mixture components, in which the three components are denoted red, blue, and green, and the values of the mixing coefficients are shown below each component. (b) Contours of the marginal probability density  $p(\mathbf{x})$  of the mixture distribution. (c) A surface plot of the distribution  $p(\mathbf{x})$ .

From the sum and product rules of probability, the marginal density can be written as

$$p(\mathbf{x}) = \sum_{k=1}^K p(k)p(\mathbf{x}|k), \quad (3.114)$$

which is equivalent to (3.111) in which we can view  $\pi_k = p(k)$  as the prior probability of picking the  $k$ th component, and the density  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = p(\mathbf{x}|k)$  as the probability of  $\mathbf{x}$  conditioned on  $k$ . As we will see in later chapters, an important role is played by the corresponding posterior probabilities  $p(k|\mathbf{x})$ , which are also known as *responsibilities*. From Bayes' theorem, these are given by

$$\begin{aligned} \gamma_k(\mathbf{x}) &\equiv p(k|\mathbf{x}) \\ &= \frac{p(k)p(\mathbf{x}|k)}{\sum_l p(l)p(\mathbf{x}|l)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_l \pi_l \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}. \end{aligned} \quad (3.115)$$

The form of the Gaussian mixture distribution is governed by the parameters  $\boldsymbol{\pi}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\Sigma}$ , where we have used the notation  $\boldsymbol{\pi} \equiv \{\pi_1, \dots, \pi_K\}$ ,  $\boldsymbol{\mu} \equiv \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\}$ , and  $\boldsymbol{\Sigma} \equiv \{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K\}$ . One way to set the values of these parameters is to use maximum likelihood. From (3.111), the log of the likelihood function is given by

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (3.116)$$

where  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ . We immediately see that the situation is now much more complex than with a single Gaussian, due to the summation over  $k$  inside the log-



arithm. As a result, the maximum likelihood solution for the parameters no longer has a closed-form analytical solution. One approach for maximizing the likelihood function is to use iterative numerical optimization techniques. Alternatively, we can employ a powerful framework called *expectation maximization*, which has wide applicability to a variety of different deep generative models.

### 3.3. Periodic Variables

Although Gaussian distributions are of great practical significance, both in their own right and as building blocks for more complex probabilistic models, there are situations in which they are inappropriate as density models for continuous variables. One important case, which arises in practical applications, is that of periodic variables.

An example of a periodic variable is the wind direction at a particular geographical location. We might, for instance, measure the wind direction at multiple locations and wish to summarize this data using a parametric distribution. Another example is calendar time, where we may be interested in modelling quantities that are believed to be periodic over 24 hours or over an annual cycle. Such quantities can conveniently be represented using an angular (polar) coordinate  $0 \leq \theta < 2\pi$ .

We might be tempted to treat periodic variables by choosing some direction as the origin and then applying a conventional distribution such as the Gaussian. Such an approach, however, would give results that were strongly dependent on the arbitrary choice of origin. Suppose, for instance, that we have two observations at  $\theta_1 = 1^\circ$  and  $\theta_2 = 359^\circ$ , and we model them using a standard univariate Gaussian distribution. If we place the origin at  $0^\circ$ , then the sample mean of this data set will be  $180^\circ$  with standard deviation  $179^\circ$ , whereas if we place the origin at  $180^\circ$ , then the mean will be  $0^\circ$  and the standard deviation will be  $1^\circ$ . We clearly need to develop a special approach for periodic variables.

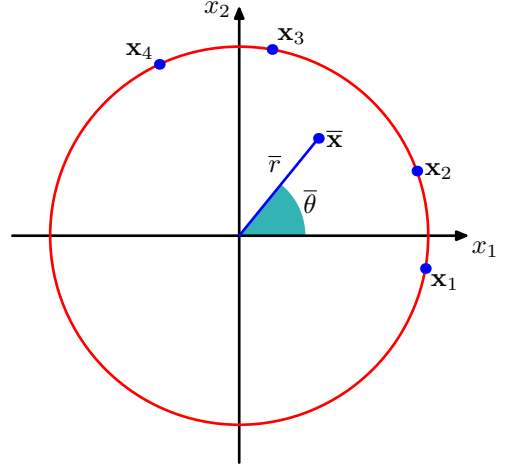
#### 3.3.1 Von Mises distribution

Let us consider the problem of evaluating the mean of a set of observations  $\mathcal{D} = \{\theta_1, \dots, \theta_N\}$  of a periodic variable  $\theta$  where  $\theta$  is measured in radians. We have already seen that the simple average  $(\theta_1 + \dots + \theta_N)/N$  will be strongly coordinate dependent. To find an invariant measure of the mean, note that the observations can be viewed as points on the unit circle and can therefore be described instead by two-dimensional unit vectors  $\mathbf{x}_1, \dots, \mathbf{x}_N$  where  $\|\mathbf{x}_n\| = 1$  for  $n = 1, \dots, N$ , as illustrated in Figure 3.9. We can average the vectors  $\{\mathbf{x}_n\}$  instead to give

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (3.117)$$

and then find the corresponding angle  $\bar{\theta}$  of this average. Clearly, this definition will ensure that the location of the mean is independent of the origin of the angular coordinate. Note that  $\bar{\mathbf{x}}$  will typically lie inside the unit circle. The Cartesian coordinates

**Figure 3.9** Illustration of the representation of values  $\theta_n$  of a periodic variable as two-dimensional vectors  $\mathbf{x}_n$  living on the unit circle. Also shown is the average  $\bar{\mathbf{x}}$  of those vectors.



of the observations are given by  $\mathbf{x}_n = (\cos \theta_n, \sin \theta_n)$ , and we can write the Cartesian coordinates of the sample mean in the form  $\bar{\mathbf{x}} = (\bar{r} \cos \bar{\theta}, \bar{r} \sin \bar{\theta})$ . Substituting into (3.117) and equating the  $x_1$  and  $x_2$  components then gives

$$\bar{x}_1 = \bar{r} \cos \bar{\theta} = \frac{1}{N} \sum_{n=1}^N \cos \theta_n, \quad \bar{x}_2 = \bar{r} \sin \bar{\theta} = \frac{1}{N} \sum_{n=1}^N \sin \theta_n. \quad (3.118)$$

Taking the ratio, and using the identity  $\tan \theta = \sin \theta / \cos \theta$ , we can solve for  $\bar{\theta}$  to give

$$\bar{\theta} = \tan^{-1} \left\{ \frac{\sum_n \sin \theta_n}{\sum_n \cos \theta_n} \right\}. \quad (3.119)$$

Shortly, we will see how this result arises naturally as a maximum likelihood estimator.

First, we need to define a periodic generalization of the Gaussian called the *von Mises* distribution. Here we will limit our attention to univariate distributions, although analogous periodic distributions can also be found over hyperspheres of arbitrary dimension (Mardia and Jupp, 2000).

By convention, we will consider distributions  $p(\theta)$  that have period  $2\pi$ . Any probability density  $p(\theta)$  defined over  $\theta$  must not only be non-negative and integrate to one, but it must also be periodic. Thus,  $p(\theta)$  must satisfy the three conditions:

$$p(\theta) \geq 0 \quad (3.120)$$

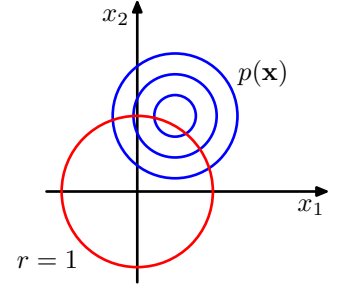
$$\int_0^{2\pi} p(\theta) d\theta = 1 \quad (3.121)$$

$$p(\theta + 2\pi) = p(\theta). \quad (3.122)$$

From (3.122), it follows that  $p(\theta + M2\pi) = p(\theta)$  for any integer  $M$ .

We can easily obtain a Gaussian-like distribution that satisfies these three properties as follows. Consider a Gaussian distribution over two variables  $\mathbf{x} = (x_1, x_2)$

**Figure 3.10** The von Mises distribution can be derived by considering a two-dimensional Gaussian of the form (3.123), whose density contours are shown in blue, and conditioning on the unit circle shown in red.



having mean  $\boldsymbol{\mu} = (\mu_1, \mu_2)$  and a covariance matrix  $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$  where  $\mathbf{I}$  is the  $2 \times 2$  identity matrix, so that

$$p(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp \left\{ -\frac{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}{2\sigma^2} \right\}. \quad (3.123)$$

The contours of constant  $p(\mathbf{x})$  are circles, as illustrated in Figure 3.10.

Now suppose we consider the value of this distribution along a circle of fixed radius. Then by construction, this distribution will be periodic, although it will not be normalized. We can determine the form of this distribution by transforming from Cartesian coordinates  $(x_1, x_2)$  to polar coordinates  $(r, \theta)$  so that

$$x_1 = r \cos \theta, \quad x_2 = r \sin \theta. \quad (3.124)$$

We also map the mean  $\boldsymbol{\mu}$  into polar coordinates by writing

$$\mu_1 = r_0 \cos \theta_0, \quad \mu_2 = r_0 \sin \theta_0. \quad (3.125)$$

Next we substitute these transformations into the two-dimensional Gaussian distribution (3.123), and then condition on the unit circle  $r = 1$ , noting that we are interested only in the dependence on  $\theta$ . Focusing on the exponent in the Gaussian distribution we have

$$\begin{aligned} & -\frac{1}{2\sigma^2} \{ (r \cos \theta - r_0 \cos \theta_0)^2 + (r \sin \theta - r_0 \sin \theta_0)^2 \} \\ &= -\frac{1}{2\sigma^2} \{ 1 + r_0^2 - 2r_0 \cos \theta \cos \theta_0 - 2r_0 \sin \theta \sin \theta_0 \} \\ &= \frac{r_0}{\sigma^2} \cos(\theta - \theta_0) + \text{const} \end{aligned} \quad (3.126)$$

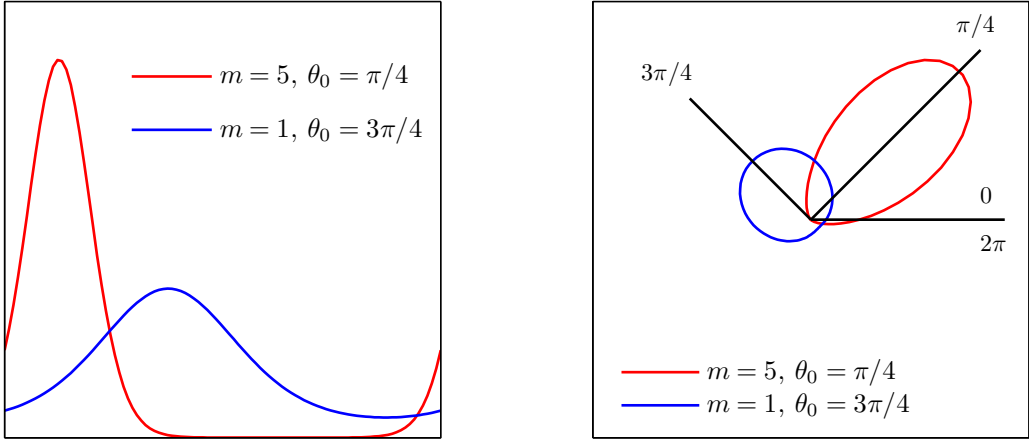
where ‘const’ denotes terms independent of  $\theta$ . We have made use of the following trigonometrical identities:

$$\cos^2 A + \sin^2 A = 1 \quad (3.127)$$

$$\cos A \cos B + \sin A \sin B = \cos(A - B). \quad (3.128)$$

If we now define  $m = r_0/\sigma^2$ , we obtain our final expression for the distribution of  $p(\theta)$  along the unit circle  $r = 1$  in the form

$$p(\theta|\theta_0, m) = \frac{1}{2\pi I_0(m)} \exp \{ m \cos(\theta - \theta_0) \}, \quad (3.129)$$



**Figure 3.11** The von Mises distribution plotted for two different parameter values, shown as a Cartesian plot on the left and as the corresponding polar plot on the right.

which is called the *von Mises* distribution or the *circular normal*. Here the parameter  $\theta_0$  corresponds to the mean of the distribution, whereas  $m$ , which is known as the *concentration* parameter, is analogous to the inverse variance (i.e. the precision) for the Gaussian. The normalization coefficient in (3.129) is expressed in terms of  $I_0(m)$ , which is the zeroth-order modified Bessel function of the first kind (Abramowitz and Stegun, 1965) and is defined by

$$I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp \{m \cos \theta\} d\theta. \quad (3.130)$$

**Exercise 3.31**

For large  $m$ , the distribution becomes approximately Gaussian. The von Mises distribution is plotted in Figure 3.11, and the function  $I_0(m)$  is plotted in Figure 3.12.

Now consider the maximum likelihood estimators for the parameters  $\theta_0$  and  $m$  for the von Mises distribution. The log likelihood function is given by

$$\ln p(\mathcal{D}|\theta_0, m) = -N \ln(2\pi) - N \ln I_0(m) + m \sum_{n=1}^N \cos(\theta_n - \theta_0). \quad (3.131)$$

Setting the derivative with respect to  $\theta_0$  equal to zero gives

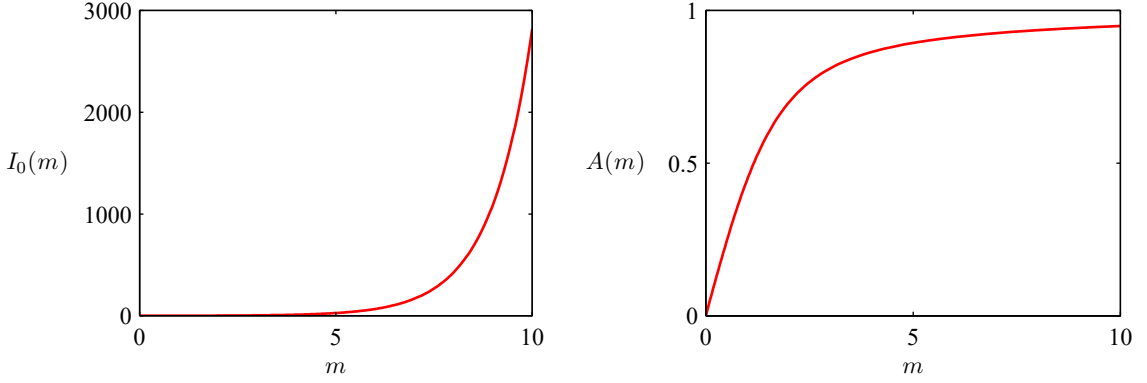
$$\sum_{n=1}^N \sin(\theta_n - \theta_0) = 0. \quad (3.132)$$

To solve for  $\theta_0$ , we make use of the trigonometric identity

$$\sin(A - B) = \cos B \sin A - \sin B \cos A \quad (3.133)$$

**Exercise 3.32**

from which we obtain



**Figure 3.12** Plot of the Bessel function  $I_0(m)$  defined by (3.130), together with the function  $A(m)$  defined by (3.136).

$$\theta_0^{\text{ML}} = \tan^{-1} \left\{ \frac{\sum_n \sin \theta_n}{\sum_n \cos \theta_n} \right\}, \quad (3.134)$$

which we recognize as the result (3.119) obtained earlier for the mean of the observations viewed in a two-dimensional Cartesian space.

Similarly, maximizing (3.131) with respect to  $m$  and making use of  $I'_0(m) = I_1(m)$  (Abramowitz and Stegun, 1965), we have

$$A(m_{\text{ML}}) = \frac{1}{N} \sum_{n=1}^N \cos(\theta_n - \theta_0^{\text{ML}}) \quad (3.135)$$

where we have substituted for the maximum likelihood solution for  $\theta_0^{\text{ML}}$  (recalling that we are performing a joint optimization over  $\theta$  and  $m$ ), and we have defined

$$A(m) = \frac{I_1(m)}{I_0(m)}. \quad (3.136)$$

The function  $A(m)$  is plotted in Figure 3.12. Making use of the trigonometric identity (3.128), we can write (3.135) in the form

$$A(m_{\text{ML}}) = \left( \frac{1}{N} \sum_{n=1}^N \cos \theta_n \right) \cos \theta_0^{\text{ML}} + \left( \frac{1}{N} \sum_{n=1}^N \sin \theta_n \right) \sin \theta_0^{\text{ML}}. \quad (3.137)$$

The right-hand side of (3.137) is easily evaluated, and the function  $A(m)$  can be inverted numerically. One limitation of the von Mises distribution is that it is unimodal. By forming *mixtures* of von Mises distributions, we obtain a flexible framework for modelling periodic variables that can handle multimodality.

For completeness, we mention briefly some alternative techniques for constructing periodic distributions. The simplest approach is to use a histogram of observations in which the angular coordinate is divided into fixed bins. This has the virtue of

## Section 3.5

simplicity and flexibility but also suffers from significant limitations, as we will see when we discuss histogram methods in more detail later. Another approach starts, like the von Mises distribution, from a Gaussian distribution over a Euclidean space but now marginalizes onto the unit circle rather than conditioning (Mardia and Jupp, 2000). However, this leads to more complex forms of distribution and will not be discussed further. Finally, any valid distribution over the real axis (such as a Gaussian) can be turned into a periodic distribution by mapping successive intervals of width  $2\pi$  onto the periodic variable  $(0, 2\pi)$ , which corresponds to ‘wrapping’ the real axis around the unit circle. Again, the resulting distribution is more complex to handle than the von Mises distribution.

### 3.4. The Exponential Family

The probability distributions that we have studied so far in this chapter (with the exception of mixture models) are specific examples of a broad class of distributions called the *exponential family* (Duda and Hart, 1973; Bernardo and Smith, 1994). Members of the exponential family have many important properties in common, and it is illuminating to discuss these properties in some generality.

The exponential family of distributions over  $\mathbf{x}$ , given parameters  $\boldsymbol{\eta}$ , is defined to be the set of distributions of the form

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} \quad (3.138)$$

where  $\mathbf{x}$  may be scalar or vector and may be discrete or continuous. Here  $\boldsymbol{\eta}$  are called the *natural parameters* of the distribution, and  $\mathbf{u}(\mathbf{x})$  is some function of  $\mathbf{x}$ . The function  $g(\boldsymbol{\eta})$  can be interpreted as the coefficient that ensures that the distribution is normalized, and therefore, it satisfies

$$g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} d\mathbf{x} = 1 \quad (3.139)$$

where the integration is replaced by summation if  $\mathbf{x}$  is a discrete variable.

We begin by taking some examples of the distributions introduced earlier in the chapter and showing that they are indeed members of the exponential family. Consider first the Bernoulli distribution:

$$p(x|\mu) = \text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}. \quad (3.140)$$

Expressing the right-hand side as the exponential of the logarithm, we have

$$\begin{aligned} p(x|\mu) &= \exp \{ x \ln \mu + (1 - x) \ln(1 - \mu) \} \\ &= (1 - \mu) \exp \left\{ \ln \left( \frac{\mu}{1 - \mu} \right) x \right\}. \end{aligned} \quad (3.141)$$

Comparison with (3.138) allows us to identify

$$\eta = \ln \left( \frac{\mu}{1 - \mu} \right) \quad (3.142)$$

which we can solve for  $\mu$  to give  $\mu = \sigma(\eta)$ , where

$$\sigma(\eta) = \frac{1}{1 + \exp(-\eta)} \quad (3.143)$$

is called the *logistic sigmoid* function. Thus, we can write the Bernoulli distribution using the standard representation (3.138) in the form

$$p(x|\eta) = \sigma(-\eta) \exp(\eta x) \quad (3.144)$$

where we have used  $1 - \sigma(\eta) = \sigma(-\eta)$ , which is easily proved from (3.143). Comparison with (3.138) shows that

$$u(x) = x \quad (3.145)$$

$$h(x) = 1 \quad (3.146)$$

$$g(\eta) = \sigma(-\eta). \quad (3.147)$$

Next consider the multinomial distribution which, for a single observation  $\mathbf{x}$ , takes the form

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^M \mu_k^{x_k} = \exp \left\{ \sum_{k=1}^M x_k \ln \mu_k \right\} \quad (3.148)$$

where  $\mathbf{x} = (x_1, \dots, x_M)^T$ . Again, we can write this in the standard representation (3.138) so that

$$p(\mathbf{x}|\boldsymbol{\eta}) = \exp(\boldsymbol{\eta}^T \mathbf{x}) \quad (3.149)$$

where  $\eta_k = \ln \mu_k$ , and we have defined  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_M)^T$ . Again, comparing with (3.138) we have

$$\mathbf{u}(\mathbf{x}) = \mathbf{x} \quad (3.150)$$

$$h(\mathbf{x}) = 1 \quad (3.151)$$

$$g(\boldsymbol{\eta}) = 1. \quad (3.152)$$

Note that the parameters  $\eta_k$  are not independent because the parameters  $\mu_k$  are subject to the constraint

$$\sum_{k=1}^M \mu_k = 1 \quad (3.153)$$

so that, given any  $M - 1$  of the parameters  $\mu_k$ , the value of the remaining parameter is fixed. In some circumstances, it will be convenient to remove this constraint by expressing the distribution in terms of only  $M - 1$  parameters. This can be achieved by using the relationship (3.153) to eliminate  $\mu_M$  by expressing it in terms of the remaining  $\{\mu_k\}$  where  $k = 1, \dots, M - 1$ , thereby leaving  $M - 1$  parameters. Note that these remaining parameters are still subject to the constraints

$$0 \leq \mu_k \leq 1, \quad \sum_{k=1}^{M-1} \mu_k \leq 1. \quad (3.154)$$

Making use of the constraint (3.153), the multinomial distribution in this representation then becomes

$$\begin{aligned}
 & \exp \left\{ \sum_{k=1}^M x_k \ln \mu_k \right\} \\
 &= \exp \left\{ \sum_{k=1}^{M-1} x_k \ln \mu_k + \left( 1 - \sum_{k=1}^{M-1} x_k \right) \ln \left( 1 - \sum_{k=1}^{M-1} \mu_k \right) \right\} \\
 &= \exp \left\{ \sum_{k=1}^{M-1} x_k \ln \left( \frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j} \right) + \ln \left( 1 - \sum_{k=1}^{M-1} \mu_k \right) \right\}. \quad (3.155)
 \end{aligned}$$

We now identify

$$\ln \left( \frac{\mu_k}{1 - \sum_j \mu_j} \right) = \eta_k, \quad (3.156)$$

which we can solve for  $\mu_k$  by first summing both sides over  $k$  and then rearranging and back-substituting to give

$$\mu_k = \frac{\exp(\eta_k)}{1 + \sum_j \exp(\eta_j)}. \quad (3.157)$$

This is called the *softmax* function or the *normalized exponential*. In this representation, the multinomial distribution therefore takes the form

$$p(\mathbf{x}|\boldsymbol{\eta}) = \left( 1 + \sum_{k=1}^{M-1} \exp(\eta_k) \right)^{-1} \exp(\boldsymbol{\eta}^T \mathbf{x}). \quad (3.158)$$

This is the standard form of the exponential family, with parameter vector  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_{M-1})^T$  in which

$$\mathbf{u}(\mathbf{x}) = \mathbf{x} \quad (3.159)$$

$$h(\mathbf{x}) = 1 \quad (3.160)$$

$$g(\boldsymbol{\eta}) = \left( 1 + \sum_{k=1}^{M-1} \exp(\eta_k) \right)^{-1}. \quad (3.161)$$

Finally, let us consider the Gaussian distribution. For the univariate Gaussian, we have

$$p(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} \quad (3.162)$$

$$= \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} x^2 + \frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} \mu^2 \right\}, \quad (3.163)$$

which, after some simple rearranging, can be cast in the standard exponential family form (3.138) with

*Exercise 3.35*



$$\boldsymbol{\eta} = \begin{pmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{pmatrix} \quad (3.164)$$

$$\mathbf{u}(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix} \quad (3.165)$$

$$h(\mathbf{x}) = (2\pi)^{-1/2} \quad (3.166)$$

$$g(\boldsymbol{\eta}) = (-2\eta_2)^{1/2} \exp\left(\frac{\eta_1^2}{4\eta_2}\right). \quad (3.167)$$

Finally, we shall sometimes make use of a restricted form of (3.138) in which we choose  $\mathbf{u}(\mathbf{x}) = \mathbf{x}$ . However, this can be somewhat generalized by noting that if  $f(\mathbf{x})$  is a normalized density then

$$\frac{1}{s} f\left(\frac{1}{s}\mathbf{x}\right) \quad (3.168)$$

is also a normalized density, where  $s > 0$  is a scale parameter. Combining these, we arrive at a restricted set of exponential family class-conditional densities of the form

$$p(\mathbf{x}|\boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s}\mathbf{x}\right) g(\boldsymbol{\lambda}_k) \exp\left\{\frac{1}{s}\boldsymbol{\lambda}_k^T \mathbf{x}\right\}. \quad (3.169)$$

Note that we are allowing each class to have its own parameter vector  $\boldsymbol{\lambda}_k$  but we are assuming that the classes share the same scale parameter  $s$ .

### 3.4.1 Sufficient statistics

Let us now consider the problem of estimating the parameter vector  $\boldsymbol{\eta}$  in the general exponential family distribution (3.138) using the technique of maximum likelihood. Taking the gradient of both sides of (3.139) with respect to  $\boldsymbol{\eta}$ , we have

$$\begin{aligned} \nabla g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})\} d\mathbf{x} \\ + g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})\} \mathbf{u}(\mathbf{x}) d\mathbf{x} = 0. \end{aligned} \quad (3.170)$$

Rearranging and making use again of (3.139) then gives

$$-\frac{1}{g(\boldsymbol{\eta})} \nabla g(\boldsymbol{\eta}) = \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})\} \mathbf{u}(\mathbf{x}) d\mathbf{x} = \mathbb{E}[\mathbf{u}(\mathbf{x})]. \quad (3.171)$$

We therefore obtain the result

$$-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}[\mathbf{u}(\mathbf{x})]. \quad (3.172)$$

Note that the covariance of  $\mathbf{u}(\mathbf{x})$  can be expressed in terms of the second derivatives of  $g(\boldsymbol{\eta})$ , and similarly for higher-order moments. Thus, provided we can normalize a distribution from the exponential family, we can always find its moments by simple differentiation.

#### Exercise 3.36

Now consider a set of independent identically distributed data denoted by  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , for which the likelihood function is given by

$$p(\mathbf{X}|\boldsymbol{\eta}) = \left( \prod_{n=1}^N h(\mathbf{x}_n) \right) g(\boldsymbol{\eta})^N \exp \left\{ \boldsymbol{\eta}^T \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) \right\}. \quad (3.173)$$

Setting the gradient of  $\ln p(\mathbf{X}|\boldsymbol{\eta})$  with respect to  $\boldsymbol{\eta}$  to zero, we get the following condition to be satisfied by the maximum likelihood estimator  $\boldsymbol{\eta}_{\text{ML}}$ :

$$-\nabla \ln g(\boldsymbol{\eta}_{\text{ML}}) = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n), \quad (3.174)$$

which can in principle be solved to obtain  $\boldsymbol{\eta}_{\text{ML}}$ . We see that the solution for the maximum likelihood estimator depends on the data only through  $\sum_n \mathbf{u}(\mathbf{x}_n)$ , which is therefore called the *sufficient statistic* of the distribution (3.138). We do not need to store the entire data set itself but only the value of the sufficient statistic. For the Bernoulli distribution, for example, the function  $\mathbf{u}(x)$  is given just by  $x$  and so we need only keep the sum of the data points  $\{x_n\}$ , whereas for the Gaussian  $\mathbf{u}(x) = (x, x^2)^T$ , and so we should keep both the sum of  $\{x_n\}$  and the sum of  $\{x_n^2\}$ .

If we consider the limit  $N \rightarrow \infty$ , then the right-hand side of (3.174) becomes  $\mathbb{E}[\mathbf{u}(\mathbf{x})]$ , and so by comparing with (3.172) we see that in this limit,  $\boldsymbol{\eta}_{\text{ML}}$  will equal the true value  $\boldsymbol{\eta}$ .

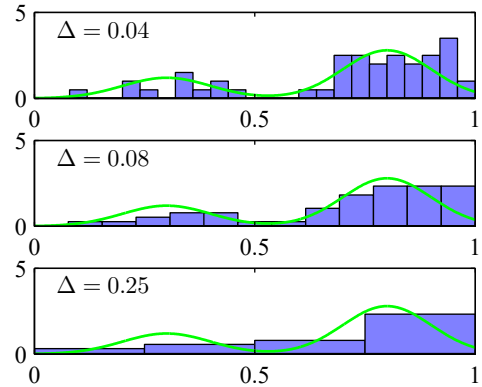
### 3.5. Nonparametric Methods

Throughout this chapter, we have focused on the use of probability distributions having specific functional forms governed by a small number of parameters whose values are to be determined from a data set. This is called the *parametric* approach to density modelling. An important limitation of this approach is that the chosen density might be a poor model of the distribution that generates the data, which can result in poor predictive performance. For instance, if the process that generates the data is multimodal, then this aspect of the distribution can never be captured by a Gaussian, which is necessarily unimodal. In this final section, we consider some *nonparametric* approaches to density estimation that make few assumptions about the form of the distribution.

#### 3.5.1 Histograms

Let us start with a discussion of histogram methods for density estimation, which we have already encountered in the context of marginal and conditional distributions in Figure 2.5 and in the context of the central limit theorem in Figure 3.2. Here we explore the properties of histogram density models in more detail, focusing on cases with a single continuous variable  $x$ . Standard histograms simply partition  $x$  into distinct bins of width  $\Delta_i$  and then count the number  $n_i$  of observations of  $x$  falling

**Figure 3.13** An illustration of the histogram approach to density estimation, in which a data set of 50 data points is generated from the distribution shown by the green curve. Histogram density estimates, based on (3.175) with a common bin width  $\Delta$ , are shown for various values of  $\Delta$ .



in bin  $i$ . To turn this count into a normalized probability density, we simply divide by the total number  $N$  of observations and by the width  $\Delta_i$  of the bins to obtain probability values for each bin:

$$p_i = \frac{n_i}{N\Delta_i} \quad (3.175)$$

for which it is easily seen that  $\int p(x) dx = 1$ . This gives a model for the density  $p(x)$  that is constant over the width of each bin. Often the bins are chosen to have the same width  $\Delta_i = \Delta$ .

In Figure 3.13, we show an example of histogram density estimation. Here the data is drawn from the distribution corresponding to the green curve, which is formed from a mixture of two Gaussians. Also shown are three examples of histogram density estimates corresponding to three different choices for the bin width  $\Delta$ . We see that when  $\Delta$  is very small (top figure), the resulting density model is very spiky, with a lot of structure that is not present in the underlying distribution that generated the data set. Conversely, if  $\Delta$  is too large (bottom figure) then the result is a model that is too smooth and consequently fails to capture the bimodal property of the green curve. The best results are obtained for some intermediate value of  $\Delta$  (middle figure). In principle, a histogram density model is also dependent on the choice of edge location for the bins, though this is typically much less significant than the bin width  $\Delta$ .

Note that the histogram method has the property (unlike the methods to be discussed shortly) that, once the histogram has been computed, the data set itself can be discarded, which can be advantageous if the data set is large. Also, the histogram approach is easily applied if the data points arrive sequentially.

In practice, the histogram technique can be useful for obtaining a quick visualization of data in one or two dimensions but is unsuited to most density estimation applications. One obvious problem is that the estimated density has discontinuities that are due to the bin edges rather than any property of the underlying distribution that generated the data. A major limitation of the histogram approach is its scaling with dimensionality. If we divide each variable in a  $D$ -dimensional space into

## Section 6.1.1

$M$  bins, then the total number of bins will be  $M^D$ . This exponential scaling with  $D$  is an example of the *curse of dimensionality*. In a space of high dimensionality, the quantity of data needed to provide meaningful estimates of the local probability density would be prohibitive.

## Chapter 1

The histogram approach to density estimation does, however, teach us two important lessons. First, to estimate the probability density at a particular location, we should consider the data points that lie within some local neighbourhood of that point. Note that the concept of locality requires that we assume some form of distance measure, and here we have been assuming Euclidean distance. For histograms, this neighbourhood property was defined by the bins, and there is a natural ‘smoothing’ parameter describing the spatial extent of the local region, in this case the bin width. Second, to obtain good results, the value of the smoothing parameter should be neither too large nor too small. This is reminiscent of the choice of model complexity in polynomial regression where the degree  $M$  of the polynomial, or alternatively the value  $\lambda$  of the regularization parameter, was optimal for some intermediate value, neither too large nor too small. Armed with these insights, we turn now to a discussion of two widely used nonparametric techniques for density estimation, kernel estimators and nearest neighbours, which have better scaling with dimensionality than the simple histogram model.

## 3.5.2 Kernel densities

Let us suppose that observations are being drawn from some unknown probability density  $p(\mathbf{x})$  in some  $D$ -dimensional space, which we will take to be Euclidean, and we wish to estimate the value of  $p(\mathbf{x})$ . From our earlier discussion of locality, let us consider some small region  $\mathcal{R}$  containing  $\mathbf{x}$ . The probability mass associated with this region is given by

$$P = \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x}. \quad (3.176)$$

## Section 3.1.2

Now suppose that we have collected a data set comprising  $N$  observations drawn from  $p(\mathbf{x})$ . Because each data point has a probability  $P$  of falling within  $\mathcal{R}$ , the total number  $K$  of points that lie inside  $\mathcal{R}$  will be distributed according to the binomial distribution:

$$\text{Bin}(K|N, P) = \frac{N!}{K!(N-K)!} P^K (1-P)^{N-K}. \quad (3.177)$$

Using (3.11), we see that the mean fraction of points falling inside the region is  $\mathbb{E}[K/N] = P$ , and similarly using (3.12), we see that the variance around this mean is  $\text{var}[K/N] = P(1-P)/N$ . For large  $N$ , this distribution will be sharply peaked around the mean and so

$$K \simeq NP. \quad (3.178)$$

If, however, we also assume that the region  $\mathcal{R}$  is sufficiently small so that the probability density  $p(\mathbf{x})$  is roughly constant over the region, then we have

$$P \simeq p(\mathbf{x})V \quad (3.179)$$

where  $V$  is the volume of  $\mathcal{R}$ . Combining (3.178) and (3.179), we obtain our density estimate in the form

$$p(\mathbf{x}) = \frac{K}{NV}. \quad (3.180)$$

Note that the validity of (3.180) depends on two contradictory assumptions, namely that the region  $\mathcal{R}$  is sufficiently small that the density is approximately constant over the region and yet sufficiently large (in relation to the value of that density) that the number  $K$  of points falling inside the region is sufficient for the binomial distribution to be sharply peaked.

We can exploit the result (3.180) in two different ways. Either we can fix  $K$  and determine the value of  $V$  from the data, which gives rise to the  $K$ -nearest-neighbour technique discussed shortly, or we can fix  $V$  and determine  $K$  from the data, giving rise to the kernel approach. It can be shown that both the  $K$ -nearest-neighbour density estimator and the kernel density estimator converge to the true probability density in the limit  $N \rightarrow \infty$  provided that  $V$  shrinks with  $N$  and that  $K$  grows with  $N$ , at an appropriate rate (Duda and Hart, 1973).

We begin by discussing the kernel method in detail. To start with we take the region  $\mathcal{R}$  to be a small hypercube centred on the point  $\mathbf{x}$  at which we wish to determine the probability density. To count the number  $K$  of points falling within this region, it is convenient to define the following function:

$$k(\mathbf{u}) = \begin{cases} 1, & |u_i| \leq 1/2, \quad i = 1, \dots, D, \\ 0, & \text{otherwise,} \end{cases} \quad (3.181)$$

which represents a unit cube centred on the origin. The function  $k(\mathbf{u})$  is an example of a *kernel function*, and in this context, it is also called a *Parzen window*. From (3.181), the quantity  $k((\mathbf{x} - \mathbf{x}_n)/h)$  will be 1 if the data point  $\mathbf{x}_n$  lies inside a cube of side  $h$  centred on  $\mathbf{x}$ , and zero otherwise. The total number of data points lying inside this cube will therefore be

$$K = \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right). \quad (3.182)$$

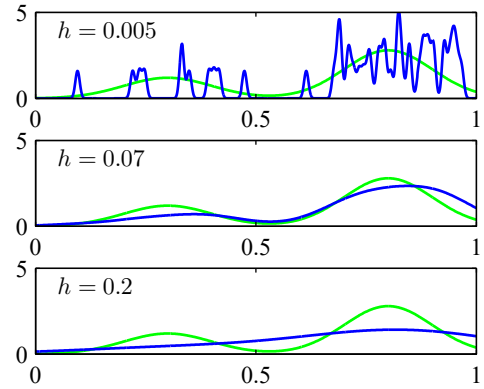
Substituting this expression into (3.180) then gives the following result for the estimated density at  $\mathbf{x}$ :

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right) \quad (3.183)$$

where we have used  $V = h^D$  for the volume of a hypercube of side  $h$  in  $D$  dimensions. Using the symmetry of the function  $k(\mathbf{u})$ , we can now reinterpret this equation, not as a single cube centred on  $\mathbf{x}$  but as the sum over  $N$  cubes centred on the  $N$  data points  $\mathbf{x}_n$ .

As it stands, the kernel density estimator (3.183) will suffer from one of the same problems that the histogram method suffered from, namely the presence of artificial discontinuities, in this case at the boundaries of the cubes. We can obtain a smoother

**Figure 3.14** Illustration of the kernel density model (3.184) applied to the same data set used to demonstrate the histogram approach in Figure 3.13. We see that  $h$  acts as a smoothing parameter and that if it is set too small (top panel), the result is a very noisy density model, whereas if it is set too large (bottom panel), then the bimodal nature of the underlying distribution from which the data is generated (shown by the green curve) is washed out. The best density model is obtained for some intermediate value of  $h$  (middle panel).



density model if we choose a smoother kernel function, and a common choice is the Gaussian, which gives rise to the following kernel density model:

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2} \right\} \quad (3.184)$$

where  $h$  represents the standard deviation of the Gaussian components. Thus, our density model is obtained by placing a Gaussian over each data point, adding up the contributions over the whole data set, and then dividing by  $N$  so that the density is correctly normalized. In Figure 3.14, we apply the model (3.184) to the data set used earlier to demonstrate the histogram technique. We see that, as expected, the parameter  $h$  plays the role of a smoothing parameter, and there is a trade-off between sensitivity to noise at small  $h$  and over-smoothing at large  $h$ . Again, the optimization of  $h$  is a problem in model complexity, analogous to the choice of bin width in histogram density estimation or the degree of the polynomial used in curve fitting.

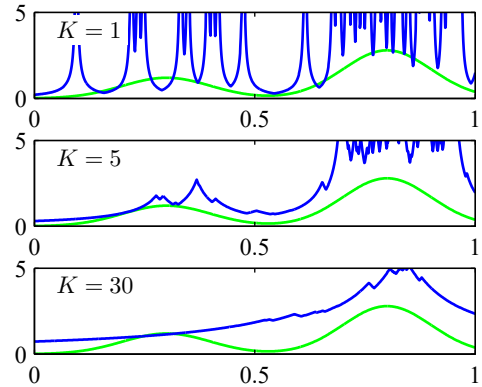
We can choose any other kernel function  $k(\mathbf{u})$  in (3.183) subject to the conditions

$$k(\mathbf{u}) \geq 0, \quad (3.185)$$

$$\int k(\mathbf{u}) d\mathbf{u} = 1, \quad (3.186)$$

which ensure that the resulting probability distribution is non-negative everywhere and integrates to one. The class of density model given by (3.183) is called a kernel density estimator or *Parzen* estimator. It has a great merit that there is no computation involved in the ‘training’ phase because this simply requires the training set to be stored. However, this is also one of its great weaknesses because the computational cost of evaluating the density grows linearly with the size of the data set.

**Figure 3.15** Illustration of  $K$ -nearest-neighbour density estimation using the same data set as in Figures 3.14 and 3.13. We see that the parameter  $K$  governs the degree of smoothing, so that a small value of  $K$  leads to a very noisy density model (top panel), whereas a large value (bottom panel) smooths out the bimodal nature of the true distribution (shown by the green curve) from which the data set was generated.



### 3.5.3 Nearest-neighbours

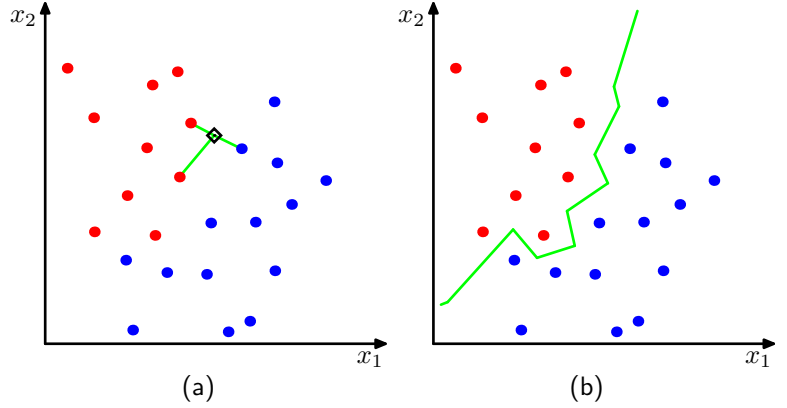
One of the difficulties with the kernel approach to density estimation is that the parameter  $h$  governing the kernel width is fixed for all kernels. In regions of high data density, a large value of  $h$  may lead to over-smoothing and a washing out of structure that might otherwise be extracted from the data. However, reducing  $h$  may lead to noisy estimates elsewhere in the data space where the density is smaller. Thus, the optimal choice for  $h$  may be dependent on the location within the data space. This issue is addressed by nearest-neighbour methods for density estimation.

We therefore return to our general result (3.180) for local density estimation, and instead of fixing  $V$  and determining the value of  $K$  from the data, we consider a fixed value of  $K$  and use the data to find an appropriate value for  $V$ . To do this, we consider a small sphere centred on the point  $\mathbf{x}$  at which we wish to estimate the density  $p(\mathbf{x})$ , and we allow the radius of the sphere to grow until it contains precisely  $K$  data points. The estimate of the density  $p(\mathbf{x})$  is then given by (3.180) with  $V$  set to the volume of the resulting sphere. This technique is known as  *$K$  nearest neighbours* and is illustrated in Figure 3.15 for various choices of the parameter  $K$  using the same data set as used in Figures 3.13 and 3.14. We see that the value of  $K$  now governs the degree of smoothing and that again there is an optimum choice for  $K$  that is neither too large nor too small. Note that the model produced by  $K$  nearest neighbours is not a true density model because the integral over all space diverges.

#### Exercise 3.38

We close this chapter by showing how the  $K$ -nearest-neighbour technique for density estimation can be extended to the problem of classification. To do this, we apply the  $K$ -nearest-neighbour density estimation technique to each class separately and then make use of Bayes' theorem. Let us suppose that we have a data set comprising  $N_k$  points in class  $\mathcal{C}_k$  with  $N$  points in total, so that  $\sum_k N_k = N$ . If we wish to classify a new point  $\mathbf{x}$ , we draw a sphere centred on  $\mathbf{x}$  containing precisely  $K$  points irrespective of their class. Suppose this sphere has volume  $V$  and contains  $K_k$  points from class  $\mathcal{C}_k$ . Then (3.180) provides an estimate of the density associated

**Figure 3.16** (a) In the  $K$ -nearest-neighbour classifier, a new point, shown by the black diamond, is classified according to the majority class membership of the  $K$  closest training data points, in this case  $K = 3$ . (b) In the nearest-neighbour ( $K = 1$ ) approach to classification, the resulting decision boundary is composed of hyperplanes that form perpendicular bisectors of pairs of points from different classes.



with each class:

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{K_k}{N_k V}. \quad (3.187)$$

Similarly, the unconditional density is given by

$$p(\mathbf{x}) = \frac{K}{NV} \quad (3.188)$$

and the class priors are given by

$$p(\mathcal{C}_k) = \frac{N_k}{N}. \quad (3.189)$$

We can now combine (3.187), (3.188), and (3.189) using Bayes' theorem to obtain the posterior probability of class membership:

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} = \frac{K_k}{K}. \quad (3.190)$$

We can minimize the probability of misclassification by assigning the test point  $\mathbf{x}$  to the class having the largest posterior probability, corresponding to the largest value of  $K_k/K$ . Thus, to classify a new point, we identify the  $K$  nearest points from the training data set and then assign the new point to the class having the largest number of representatives amongst this set. Ties can be broken at random. The particular case of  $K = 1$  is called the *nearest-neighbour* rule, because a test point is simply assigned to the same class as the nearest point from the training set. These concepts are illustrated in Figure 3.16.

An interesting property of the nearest-neighbour ( $K = 1$ ) classifier is that, in the limit  $N \rightarrow \infty$ , the error rate is never more than twice the minimum achievable error rate of an optimal classifier, i.e., one that uses the true class distributions (Cover and Hart, 1967).

As discussed so far, both the  $K$ -nearest-neighbour method and the kernel density estimator require the entire training data set to be stored, leading to expensive



computation if the data set is large. This effect can be offset, at the expense of some additional one-off computation, by constructing tree-based search structures to allow (approximate) near neighbours to be found efficiently without doing an exhaustive search of the data set. Nevertheless, these nonparametric methods are still severely limited. On the other hand, we have seen that simple parametric models are very restricted in terms of the forms of distribution that they can represent. We therefore need to find density models that are very flexible and yet for which the complexity of the models can be controlled independently of the size of the training set, and this can be achieved using deep neural networks.

## Exercises

**3.1** (★) Verify that the Bernoulli distribution (3.2) satisfies the following properties:

$$\sum_{x=0}^1 p(x|\mu) = 1 \quad (3.191)$$

$$\mathbb{E}[x] = \mu \quad (3.192)$$

$$\text{var}[x] = \mu(1 - \mu). \quad (3.193)$$

Show that the entropy  $H[x]$  of a Bernoulli-distributed random binary variable  $x$  is given by

$$H[x] = -\mu \ln \mu - (1 - \mu) \ln(1 - \mu). \quad (3.194)$$

**3.2** (★★) The form of the Bernoulli distribution given by (3.2) is not symmetric between the two values of  $x$ . In some situations, it will be more convenient to use an equivalent formulation for which  $x \in \{-1, 1\}$ , in which case the distribution can be written

$$p(x|\mu) = \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} \quad (3.195)$$

where  $\mu \in [-1, 1]$ . Show that the distribution (3.195) is normalized, and evaluate its mean, variance, and entropy.

**3.3** (★★) In this exercise, we prove that the binomial distribution (3.9) is normalized. First, use the definition (3.10) of the number of combinations of  $m$  identical objects chosen from a total of  $N$  to show that

$$\binom{N}{m} + \binom{N}{m-1} = \binom{N+1}{m}. \quad (3.196)$$

Use this result to prove by induction the following result:

$$(1+x)^N = \sum_{m=0}^N \binom{N}{m} x^m, \quad (3.197)$$