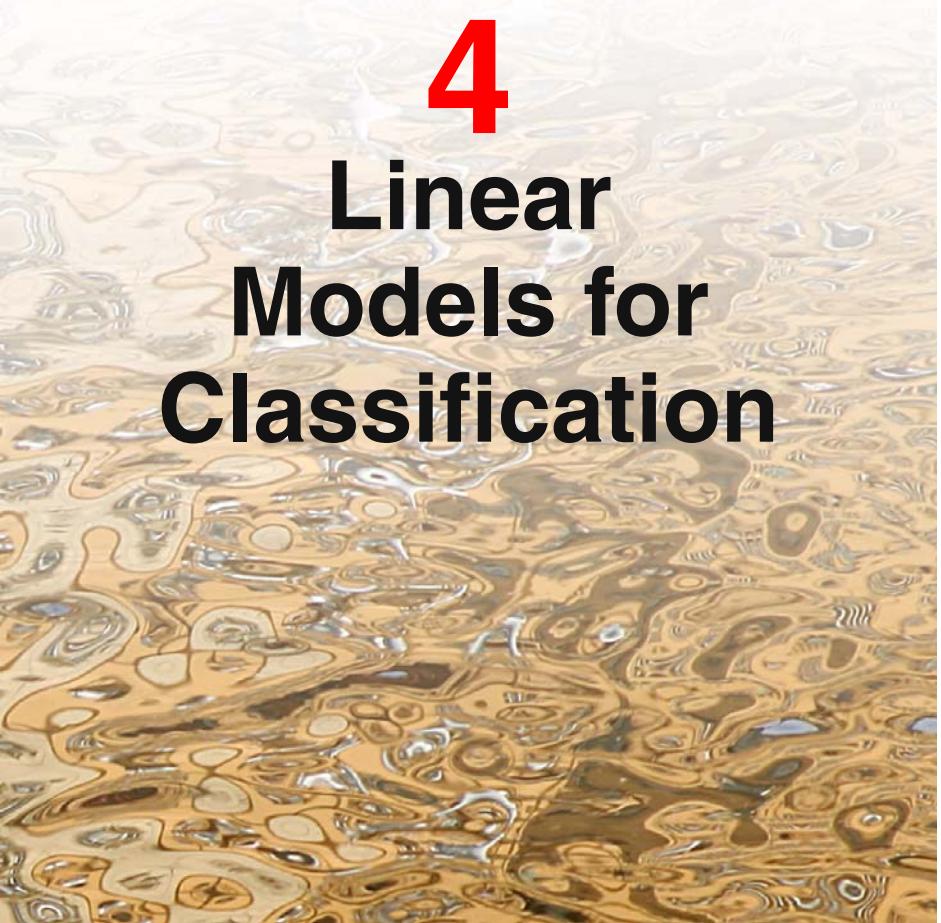


---



# 4

# Linear

# Models for

# Classification

In the previous chapter, we explored a class of regression models having particularly simple analytical and computational properties. We now discuss an analogous class of models for solving classification problems. The goal in classification is to take an input vector  $\mathbf{x}$  and to assign it to one of  $K$  discrete classes  $\mathcal{C}_k$  where  $k = 1, \dots, K$ . In the most common scenario, the classes are taken to be disjoint, so that each input is assigned to one and only one class. The input space is thereby divided into *decision regions* whose boundaries are called *decision boundaries* or *decision surfaces*. In this chapter, we consider linear models for classification, by which we mean that the decision surfaces are linear functions of the input vector  $\mathbf{x}$  and hence are defined by  $(D - 1)$ -dimensional hyperplanes within the  $D$ -dimensional input space. Data sets whose classes can be separated exactly by linear decision surfaces are said to be *linearly separable*.

For regression problems, the target variable  $t$  was simply the vector of real numbers whose values we wish to predict. In the case of classification, there are various

ways of using target values to represent class labels. For probabilistic models, the most convenient, in the case of two-class problems, is the binary representation in which there is a single target variable  $t \in \{0, 1\}$  such that  $t = 1$  represents class  $\mathcal{C}_1$  and  $t = 0$  represents class  $\mathcal{C}_2$ . We can interpret the value of  $t$  as the probability that the class is  $\mathcal{C}_1$ , with the values of probability taking only the extreme values of 0 and 1. For  $K > 2$  classes, it is convenient to use a 1-of- $K$  coding scheme in which  $\mathbf{t}$  is a vector of length  $K$  such that if the class is  $\mathcal{C}_j$ , then all elements  $t_k$  of  $\mathbf{t}$  are zero except element  $t_j$ , which takes the value 1. For instance, if we have  $K = 5$  classes, then a pattern from class 2 would be given the target vector

$$\mathbf{t} = (0, 1, 0, 0, 0)^T. \quad (4.1)$$

Again, we can interpret the value of  $t_k$  as the probability that the class is  $\mathcal{C}_k$ . For nonprobabilistic models, alternative choices of target variable representation will sometimes prove convenient.

In Chapter 1, we identified three distinct approaches to the classification problem. The simplest involves constructing a *discriminant function* that directly assigns each vector  $\mathbf{x}$  to a specific class. A more powerful approach, however, models the conditional probability distribution  $p(\mathcal{C}_k|\mathbf{x})$  in an inference stage, and then subsequently uses this distribution to make optimal decisions. By separating inference and decision, we gain numerous benefits, as discussed in Section 1.5.4. There are two different approaches to determining the conditional probabilities  $p(\mathcal{C}_k|\mathbf{x})$ . One technique is to model them directly, for example by representing them as parametric models and then optimizing the parameters using a training set. Alternatively, we can adopt a generative approach in which we model the class-conditional densities given by  $p(\mathbf{x}|\mathcal{C}_k)$ , together with the prior probabilities  $p(\mathcal{C}_k)$  for the classes, and then we compute the required posterior probabilities using Bayes' theorem

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}. \quad (4.2)$$

We shall discuss examples of all three approaches in this chapter. 3

In the linear regression models considered in Chapter 3, the model prediction  $y(\mathbf{x}, \mathbf{w})$  was given by a linear function of the parameters  $\mathbf{w}$ . In the simplest case, the model is also linear in the input variables and therefore takes the form  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ , so that  $y$  is a real number. For classification problems, however, we wish to predict discrete class labels, or more generally posterior probabilities that lie in the range  $(0, 1)$ . To achieve this, we consider a generalization of this model in which we transform the linear function of  $\mathbf{w}$  using a nonlinear function  $f(\cdot)$  so that

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0). \quad (4.3)$$

In the machine learning literature  $f(\cdot)$  is known as an *activation function*, whereas its inverse is called a *link function* in the statistics literature. The decision surfaces correspond to  $y(\mathbf{x}) = \text{constant}$ , so that  $\mathbf{w}^T \mathbf{x} + w_0 = \text{constant}$  and hence the decision surfaces are linear functions of  $\mathbf{x}$ , even if the function  $f(\cdot)$  is nonlinear. For this reason, the class of models described by (4.3) are called *generalized linear models*

(McCullagh and Nelder, 1989). Note, however, that in contrast to the models used for regression, they are no longer linear in the parameters due to the presence of the nonlinear function  $f(\cdot)$ . This will lead to more complex analytical and computational properties than for linear regression models. Nevertheless, these models are still relatively simple compared to the more general nonlinear models that will be studied in subsequent chapters.

The algorithms discussed in this chapter will be equally applicable if we first make a fixed nonlinear transformation of the input variables using a vector of basis functions  $\phi(\mathbf{x})$  as we did for regression models in Chapter 3. We begin by considering classification directly in the original input space  $\mathbf{x}$ , while in Section 4.3 we shall find it convenient to switch to a notation involving basis functions for consistency with later chapters.

## 4.1. Discriminant Functions

---

A discriminant is a function that takes an input vector  $\mathbf{x}$  and assigns it to one of  $K$  classes, denoted  $\mathcal{C}_k$ . In this chapter, we shall restrict attention to *linear discriminants*, namely those for which the decision surfaces are hyperplanes. To simplify the discussion, we consider first the case of two classes and then investigate the extension to  $K > 2$  classes.

### 4.1.1 Two classes

The simplest representation of a linear discriminant function is obtained by taking a linear function of the input vector so that

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (4.4)$$

where  $\mathbf{w}$  is called a *weight vector*, and  $w_0$  is a *bias* (not to be confused with bias in the statistical sense). The negative of the bias is sometimes called a *threshold*. An input vector  $\mathbf{x}$  is assigned to class  $\mathcal{C}_1$  if  $y(\mathbf{x}) \geq 0$  and to class  $\mathcal{C}_2$  otherwise. The corresponding decision boundary is therefore defined by the relation  $y(\mathbf{x}) = 0$ , which corresponds to a  $(D - 1)$ -dimensional hyperplane within the  $D$ -dimensional input space. Consider two points  $\mathbf{x}_A$  and  $\mathbf{x}_B$  both of which lie on the decision surface. Because  $y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0$ , we have  $\mathbf{w}^T(\mathbf{x}_A - \mathbf{x}_B) = 0$  and hence the vector  $\mathbf{w}$  is orthogonal to every vector lying within the decision surface, and so  $\mathbf{w}$  determines the orientation of the decision surface. Similarly, if  $\mathbf{x}$  is a point on the decision surface, then  $y(\mathbf{x}) = 0$ , and so the normal distance from the origin to the decision surface is given by

$$\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}. \quad (4.5)$$

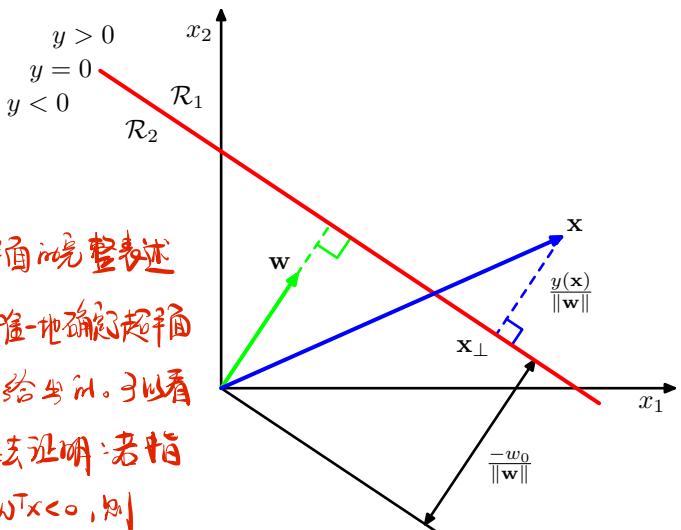
We therefore see that the bias parameter  $w_0$  determines the location of the decision surface. These properties are illustrated for the case of  $D = 2$  in Figure 4.1.

Furthermore, we note that the value of  $y(\mathbf{x})$  gives a signed measure of the perpendicular distance  $r$  of the point  $\mathbf{x}$  from the decision surface. To see this, consider

②原点所在一侧点的 $y < 0$ : 原点到超平面的带符号距离:  $\frac{w_0 + w_0}{\|w\|} = \frac{w_0}{\|w\|} < 0$

**Figure 4.1** Illustration of the geometry of a linear discriminant function in two dimensions. The decision surface, shown in red, is perpendicular to  $w$ , and its displacement from the origin is controlled by the bias parameter  $w_0$ . Also, the signed orthogonal distance of a general point  $x$  from the decision surface is given by  $y(x)/\|w\|$ .

文中未显示给定的条件是  $w_0 \leq 0$ , 即分类超平面的完整表述  
 应为  $\begin{cases} w^T x + w_0 = 0 \\ w_0 \leq 0 \end{cases}$ , 这样就剔除了正负的干扰, 唯一地确定了超平面  
 的表达式, 文中的各结论也是在  $w_0 \leq 0$  这一条件下给出的。予以看  
 到: ①由向量由原点垂直指向超平面, 及证法证明: 若指  
 向背离超平面方向, 则对超平面上任意一点  $x$ ,  $w^T x < 0$ , 则  
 $w^T x + w_0 < 0$ ,  $x$  不在超平面上的点;



an arbitrary point  $x$  and let  $x_\perp$  be its orthogonal projection onto the decision surface, so that

$$x = x_\perp + r \frac{w}{\|w\|}. \quad (4.6)$$

Multiplying both sides of this result by  $w^T$  and adding  $w_0$ , and making use of  $y(x) = w^T x + w_0$  and  $y(x_\perp) = w^T x_\perp + w_0 = 0$ , we have

$$r = \frac{y(x)}{\|w\|}. \quad (4.7)$$

This result is illustrated in Figure 4.1.

As with the linear regression models in Chapter 3, it is sometimes convenient to use a more compact notation in which we introduce an additional dummy ‘input’ value  $x_0 = 1$  and then define  $\tilde{w} = (w_0, w)$  and  $\tilde{x} = (x_0, x)$  so that

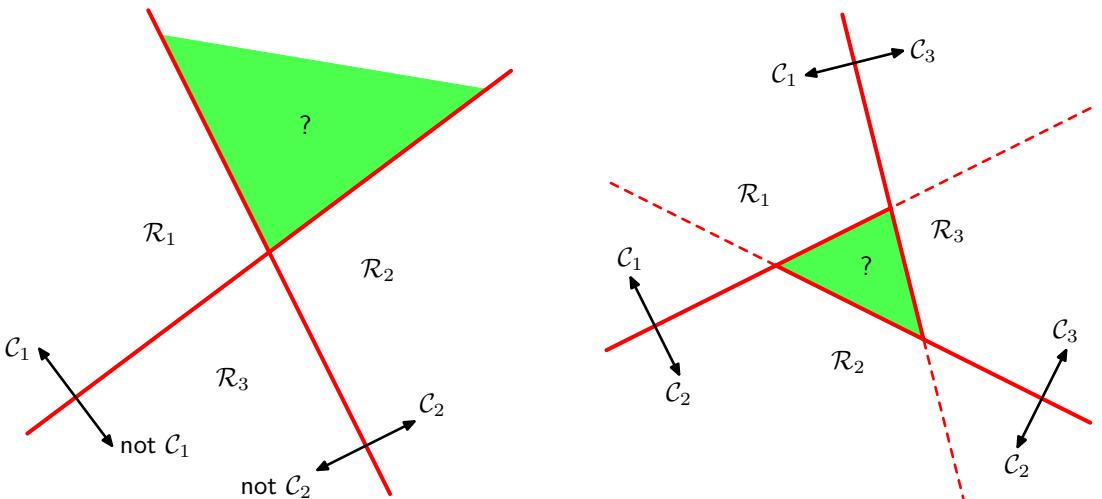
$$y(x) = \tilde{w}^T \tilde{x}. \quad (4.8)$$

In this case, the decision surfaces are  $D$ -dimensional hyperplanes passing through the origin of the  $D + 1$ -dimensional expanded input space.

### 4.1.2 Multiple classes

Now consider the extension of linear discriminants to  $K > 2$  classes. We might be tempted to build a  $K$ -class discriminant by combining a number of two-class discriminant functions. However, this leads to some serious difficulties (Duda and Hart, 1973) as we now show.

Consider the use of  $K - 1$  classifiers each of which solves a two-class problem of separating points in a particular class  $C_k$  from points not in that class. This is known as a *one-versus-the-rest* classifier. The left-hand example in Figure 4.2 shows an



**Figure 4.2** Attempting to construct a  $K$  class discriminant from a set of two class discriminants leads to ambiguous regions, shown in green. On the left is an example involving the use of two discriminants designed to distinguish points in class  $\mathcal{C}_k$  from points not in class  $\mathcal{C}_k$ . On the right is an example involving three discriminant functions each of which is used to separate a pair of classes  $\mathcal{C}_k$  and  $\mathcal{C}_j$ .

example involving three classes where this approach leads to regions of input space that are ambiguously classified.

An alternative is to introduce  $K(K - 1)/2$  binary discriminant functions, one for every possible pair of classes. This is known as a *one-versus-one* classifier. Each point is then classified according to a majority vote amongst the discriminant functions. However, this too runs into the problem of ambiguous regions, as illustrated in the right-hand diagram of Figure 4.2.

We can avoid these difficulties by considering a single  $K$ -class discriminant comprising  $K$  linear functions of the form

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (4.9)$$

and then assigning a point  $\mathbf{x}$  to class  $\mathcal{C}_k$  if  $y_k(\mathbf{x}) > y_j(\mathbf{x})$  for all  $j \neq k$ . The decision boundary between class  $\mathcal{C}_k$  and class  $\mathcal{C}_j$  is therefore given by  $y_k(\mathbf{x}) = y_j(\mathbf{x})$  and hence corresponds to a  $(D - 1)$ -dimensional hyperplane defined by

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0. \quad (4.10)$$

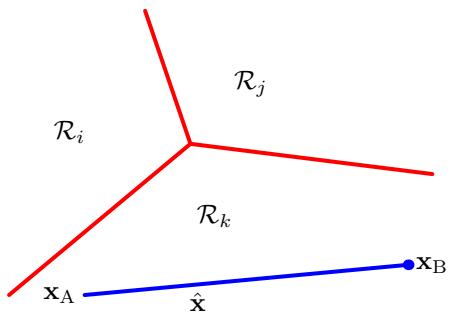
This has the same form as the decision boundary for the two-class case discussed in Section 4.1.1, and so analogous geometrical properties apply.<sup>3</sup>

The decision regions of such a discriminant are always singly connected and convex. To see this, consider two points  $\mathbf{x}_A$  and  $\mathbf{x}_B$  both of which lie inside decision region  $\mathcal{R}_k$ , as illustrated in Figure 4.3. Any point  $\hat{\mathbf{x}}$  that lies on the line connecting  $\mathbf{x}_A$  and  $\mathbf{x}_B$  can be expressed in the form

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B \quad (4.11)$$

决策区域都是  
单连通和凸的

**Figure 4.3** Illustration of the decision regions for a multiclass linear discriminant, with the decision boundaries shown in red. If two points  $\mathbf{x}_A$  and  $\mathbf{x}_B$  both lie inside the same decision region  $\mathcal{R}_k$ , then any point  $\hat{\mathbf{x}}$  that lies on the line connecting these two points must also lie in  $\mathcal{R}_k$ , and hence the decision region must be singly connected and convex.



where  $0 \leq \lambda \leq 1$ . From the linearity of the discriminant functions, it follows that

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda)y_k(\mathbf{x}_B). \quad (4.12)$$

Because both  $\mathbf{x}_A$  and  $\mathbf{x}_B$  lie inside  $\mathcal{R}_k$ , it follows that  $y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A)$ , and  $y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B)$ , for all  $j \neq k$ , and hence  $y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}})$ , and so  $\hat{\mathbf{x}}$  also lies inside  $\mathcal{R}_k$ . Thus  $\mathcal{R}_k$  is singly connected and convex.<sup>3</sup>

二分类问题也可套用这里  
多类问题的框架，其  
上小节使用单一判别器  
 $y(\mathbf{x})$ 是等价的。

Note that for two classes, we can either employ the formalism discussed here, based on two discriminant functions  $y_1(\mathbf{x})$  and  $y_2(\mathbf{x})$ , or else use the simpler but equivalent formulation described in Section 4.1.1 based on a single discriminant function  $y(\mathbf{x})$ .

We now explore three approaches to learning the parameters of linear discriminant functions, based on least squares, Fisher's linear discriminant, and the perceptron algorithm.

### 4.1.3 Least squares for classification 补充笔记已markdown笔记

In Chapter 3, we considered models that were linear functions of the parameters, and we saw that the minimization of a sum-of-squares error function led to a simple closed-form solution for the parameter values. It is therefore tempting to see if we can apply the same formalism to classification problems. Consider a general classification problem with  $K$  classes, with a 1-of- $K$  binary coding scheme for the target vector  $\mathbf{t}$ . One justification for using least squares in such a context is that it approximates the conditional expectation  $\mathbb{E}[\mathbf{t}|\mathbf{x}]$  of the target values given the input vector. For the binary coding scheme, this conditional expectation is given by the vector of posterior class probabilities. Unfortunately, however, these probabilities are typically approximated rather poorly, indeed the approximations can have values outside the range  $(0, 1)$ , due to the limited flexibility of a linear model as we shall see shortly.]

模型  $y(\mathbf{x})$  [ Each class  $\mathcal{C}_k$  is described by its own linear model so that

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (4.13)$$

where  $k = 1, \dots, K$ . We can conveniently group these together using vector notation so that

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} \quad (4.14)$$

(p+1) × K

where  $\tilde{\mathbf{W}}$  is a matrix whose  $k^{\text{th}}$  column comprises the  $D + 1$ -dimensional vector  $\tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T$  and  $\tilde{\mathbf{x}}$  is the corresponding augmented input vector  $(1, \mathbf{x}^T)^T$  with a dummy input  $x_0 = 1$ . This representation was discussed in detail in Section 3.1. A new input  $\mathbf{x}$  is then assigned to the class for which the output  $y_k = \tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}}$  is largest.

参数学习

[ We now determine the parameter matrix  $\tilde{\mathbf{W}}$  by minimizing a sum-of-squares error function, as we did for regression in Chapter 3. Consider a training data set  $\{\mathbf{x}_n, \mathbf{t}_n\}$  where  $n = 1, \dots, N$ , and define a matrix  $\mathbf{T}$  whose  $n^{\text{th}}$  row is the vector  $\mathbf{t}_n^T$ , together with a matrix  $\tilde{\mathbf{X}}$  whose  $n^{\text{th}}$  row is  $\tilde{\mathbf{x}}_n^T$ . The sum-of-squares error function can then be written as ]

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T}) \right\}. \quad (4.15)$$

Setting the derivative with respect to  $\tilde{\mathbf{W}}$  to zero, and rearranging, we then obtain the solution for  $\tilde{\mathbf{W}}$  in the form

$$\tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T} \quad (4.16)$$

where  $\tilde{\mathbf{X}}^\dagger$  is the pseudo-inverse of the matrix  $\tilde{\mathbf{X}}$ , as discussed in Section 3.1.1. We then obtain the discriminant function in the form

$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T (\tilde{\mathbf{X}}^\dagger)^T \tilde{\mathbf{x}}. \quad (4.17)$$

y(x)不满足成为一个  
概率分布的条件

[ An interesting property of least-squares solutions with multiple target variables is that if every target vector in the training set satisfies some linear constraint ]

$$\mathbf{a}^T \mathbf{t}_n + b = 0 \quad (4.18)$$

### Exercise 4.2

for some constants  $\mathbf{a}$  and  $b$ , then the model prediction for any value of  $\mathbf{x}$  will satisfy the same constraint so that

$$\mathbf{a}^T \mathbf{y}(\mathbf{x}) + b = 0. \quad (4.19)$$

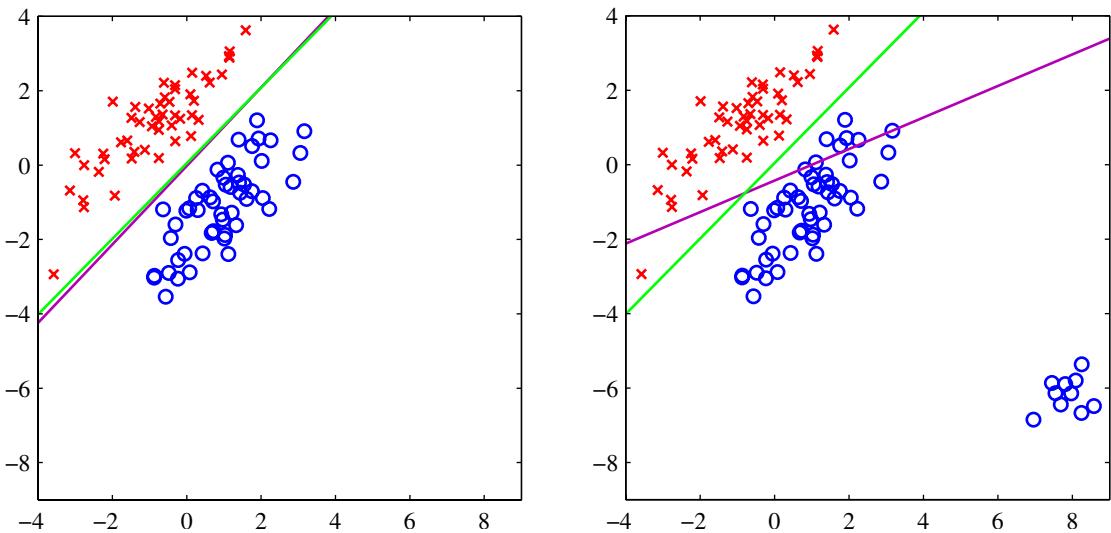
Thus if we use a 1-of- $K$  coding scheme for  $K$  classes, then the predictions made by the model will have the property that the elements of  $\mathbf{y}(\mathbf{x})$  will sum to 1 for any value of  $\mathbf{x}$ . However, this summation constraint alone is not sufficient to allow the model outputs to be interpreted as probabilities because they are not constrained to lie within the interval  $(0, 1)$ .

least-squares approach的问题

[ The least-squares approach gives an exact closed-form solution for the discriminant function parameters. However, even as a discriminant function (where we use it to make decisions directly and dispense with any probabilistic interpretation) it suffers from some severe problems. We have already seen that least-squares solutions lack robustness to outliers, and this applies equally to the classification application, as illustrated in Figure 4.4. Here we see that the additional data points in the right-hand figure produce a significant change in the location of the decision boundary, even though these points would be correctly classified by the original decision boundary in the left-hand figure. The sum-of-squares error function penalizes predictions that are ‘too correct’ in that they lie a long way on the correct side of the decision ]

### Section 2.3.7

①不能健壮对  
outliers十分敏感



**Figure 4.4** The left plot shows data from two classes, denoted by red crosses and blue circles, together with the decision boundary found by least squares (magenta curve) and also by the logistic regression model (green curve), which is discussed later in Section 4.3.2. The right-hand plot shows the corresponding results obtained when extra data points are added at the bottom left of the diagram, showing that least squares is highly sensitive to outliers, unlike logistic regression.

boundary. In Section 7.1.2, we shall consider several alternative error functions for classification and we shall see that they do not suffer from this difficulty.

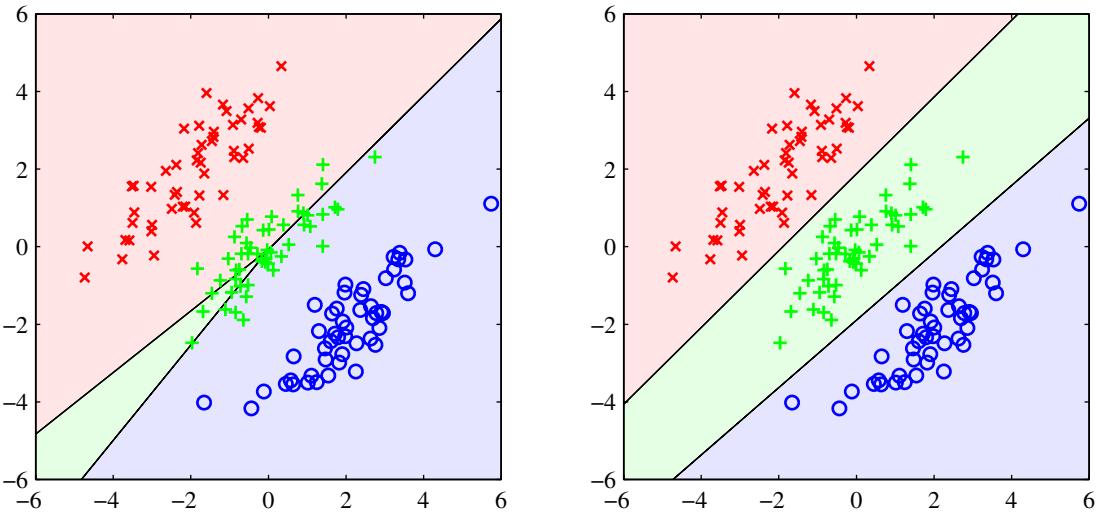
However, problems with least squares can be more severe than simply lack of robustness, as illustrated in Figure 4.5. This shows a synthetic data set drawn from three classes in a two-dimensional input space ( $x_1, x_2$ ), having the property that linear decision boundaries can give excellent separation between the classes. Indeed, the technique of logistic regression, described later in this chapter, gives a satisfactory solution as seen in the right-hand plot. However, the least-squares solution gives poor results, with only a small region of the input space assigned to the green class.

The failure of least squares should not surprise us when we recall that it corresponds to maximum likelihood under the assumption of a Gaussian conditional distribution, whereas binary target vectors clearly have a distribution that is far from Gaussian. By adopting more appropriate probabilistic models, we shall obtain classification techniques with much better properties than least squares. For the moment, however, we continue to explore alternative nonprobabilistic methods for setting the parameters in the linear classification models.] discriminants 是一种非概率方法，但由上文可见，这并不妨碍我们从概率的角度去评估、理解这些方法。

#### 4.1.4 Fisher's linear discriminant

One way to view a linear classification model is in terms of dimensionality reduction. Consider first the case of two classes, and suppose we take the  $D$ -

②模型表达能力不够，  
对有些线性可分的问题都  
无法得到满意的结果



**Figure 4.5** Example of a synthetic data set comprising three classes, with training data points denoted in red ( $\times$ ), green (+), and blue ( $\circ$ ). Lines denote the decision boundaries, and the background colours denote the respective classes of the decision regions. On the left is the result of using a least-squares discriminant. We see that the region of input space assigned to the green class is too small and so most of the points from this class are misclassified. On the right is the result of using logistic regressions as described in Section 4.3.2 showing correct classification of the training data.

dimensional input vector  $\mathbf{x}$  and project it down to one dimension using

$$y = \mathbf{w}^T \mathbf{x}. \quad (4.20)$$

If we place a threshold on  $y$  and classify  $y \geq -w_0$  as class  $\mathcal{C}_1$ , and otherwise class  $\mathcal{C}_2$ , then we obtain our standard linear classifier discussed in the previous section. In general, the projection onto one dimension leads to a considerable loss of information, and classes that are well separated in the original  $D$ -dimensional space may become strongly overlapping in one dimension. However, by adjusting the components of the weight vector  $\mathbf{w}$ , we can select a projection that maximizes the class separation. To begin with, consider a two-class problem in which there are  $N_1$  points of class  $\mathcal{C}_1$  and  $N_2$  points of class  $\mathcal{C}_2$ , so that the mean vectors of the two classes are given by

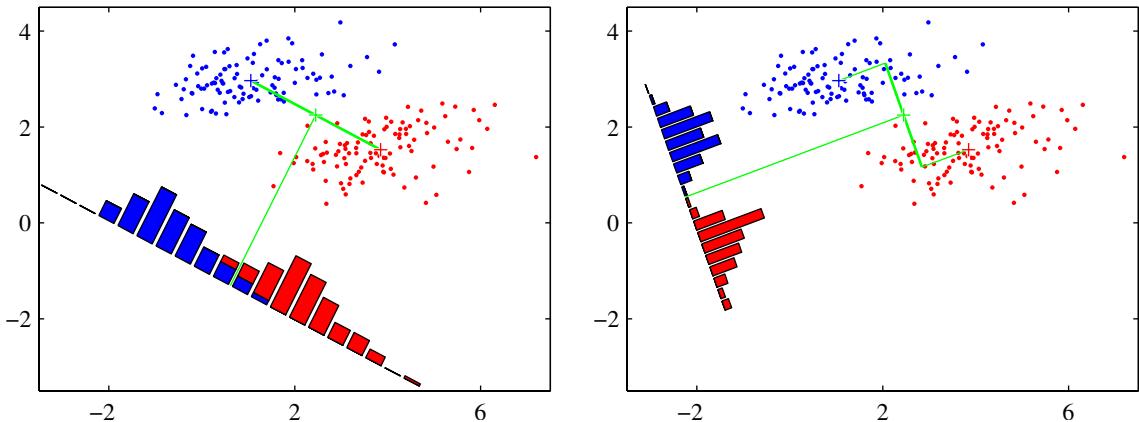
$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n. \quad (4.21)$$

The simplest measure of the separation of the classes, when projected onto  $\mathbf{w}$ , is the separation of the projected class means. This suggests that we might choose  $\mathbf{w}$  so as to maximize

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) \quad (4.22)$$

where

$$m_k = \mathbf{w}^T \mathbf{m}_k \quad (4.23)$$



**Figure 4.6** The left plot shows samples from two classes (depicted in red and blue) along with the histograms resulting from projection onto the line joining the class means. Note that there is considerable class overlap in the projected space. The right plot shows the corresponding projection based on the Fisher linear discriminant, showing the greatly improved class separation.

*Appendix E*  
*Exercise 4.4*

is the mean of the projected data from class  $\mathcal{C}_k$ . However, this expression can be made arbitrarily large simply by increasing the magnitude of  $\mathbf{w}$ . To solve this problem, we could constrain  $\mathbf{w}$  to have unit length, so that  $\sum_i w_i^2 = 1$ . Using a Lagrange multiplier to perform the constrained maximization, we then find that  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ . There is still a problem with this approach, however, as illustrated in Figure 4.6. This shows two classes that are well separated in the original two-dimensional space  $(x_1, x_2)$  but that have considerable overlap when projected onto the line joining their means. This difficulty arises from the strongly nondiagonal covariances of the class distributions. The idea proposed by Fisher is to maximize a function that will give a large separation between the projected class means while also giving a small variance within each class, thereby minimizing the class overlap.

The projection formula (4.20) transforms the set of labelled data points in  $\mathbf{x}$  into a labelled set in the one-dimensional space  $y$ . The within-class variance of the transformed data from class  $\mathcal{C}_k$  is therefore given by

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2 \quad (4.24)$$

where  $y_n = \mathbf{w}^T \mathbf{x}_n$ . We can define the total within-class variance for the whole data set to be simply  $s_1^2 + s_2^2$ . The Fisher criterion is defined to be the ratio of the between-class variance to the within-class variance and is given by

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}. \quad (4.25)$$

We can make the dependence on  $\mathbf{w}$  explicit by using (4.20), (4.23), and (4.24) to rewrite the Fisher criterion in the form

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (4.26)$$

where  $\mathbf{S}_B$  is the *between-class* covariance matrix and is given by

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \quad (4.27)$$

and  $\mathbf{S}_W$  is the total *within-class* covariance matrix, given by

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T. \quad (4.28)$$

Differentiating (4.26) with respect to  $\mathbf{w}$ , we find that  $J(\mathbf{w})$  is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}. \quad (4.29)$$

From (4.27), we see that  $\mathbf{S}_B \mathbf{w}$  is always in the direction of  $(\mathbf{m}_2 - \mathbf{m}_1)$ . Furthermore, we do not care about the magnitude of  $\mathbf{w}$ , only its direction, and so we can drop the scalar factors  $(\mathbf{w}^T \mathbf{S}_B \mathbf{w})$  and  $(\mathbf{w}^T \mathbf{S}_W \mathbf{w})$ . Multiplying both sides of (4.29) by  $\mathbf{S}_W^{-1}$  we then obtain

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1). \quad (4.30)$$

Note that if the within-class covariance is isotropic, so that  $\mathbf{S}_W$  is proportional to the unit matrix, we find that  $\mathbf{w}$  is proportional to the difference of the class means, as discussed above.

The result (4.30) is known as *Fisher's linear discriminant*, although strictly it is not a discriminant but rather a specific choice of direction for projection of the data down to one dimension. However, the projected data can subsequently be used to construct a discriminant, by choosing a threshold  $y_0$  so that we classify a new point as belonging to  $\mathcal{C}_1$  if  $y(\mathbf{x}) \geq y_0$  and classify it as belonging to  $\mathcal{C}_2$  otherwise. For example, we can model the class-conditional densities  $p(y|\mathcal{C}_k)$  using Gaussian distributions and then use the techniques of Section 1.2.4 to find the parameters of the Gaussian distributions by maximum likelihood. Having found Gaussian approximations to the projected classes, the formalism of Section 1.5.1 then gives an expression for the optimal threshold. Some justification for the Gaussian assumption comes from the central limit theorem by noting that  $y = \mathbf{w}^T \mathbf{x}$  is the sum of a set of random variables.

寻找最优 threshold  
y<sub>0</sub> 为分类点

#### 4.1.5 Relation to least squares

The least-squares approach to the determination of a linear discriminant was based on the goal of making the model predictions as close as possible to a set of target values. By contrast, the Fisher criterion was derived by requiring maximum class separation in the output space. It is interesting to see the relationship between these two approaches. In particular, we shall show that, for the two-class problem, the Fisher criterion can be obtained as a special case of least squares.

So far we have considered 1-of- $K$  coding for the target values. If, however, we adopt a slightly different target coding scheme, then the least-squares solution for

the weights becomes equivalent to the Fisher solution (Duda and Hart, 1973). In particular, we shall take the targets for class  $\mathcal{C}_1$  to be  $N/N_1$ , where  $N_1$  is the number of patterns in class  $\mathcal{C}_1$ , and  $N$  is the total number of patterns. This target value approximates the reciprocal of the prior probability for class  $\mathcal{C}_1$ . For class  $\mathcal{C}_2$ , we shall take the targets to be  $-N/N_2$ , where  $N_2$  is the number of patterns in class  $\mathcal{C}_2$ .

The sum-of-squares error function can be written

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2. \quad (4.31)$$

Setting the derivatives of  $E$  with respect to  $w_0$  and  $\mathbf{w}$  to zero, we obtain respectively

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = 0 \quad (4.32)$$

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0. \quad (4.33)$$

From (4.32), and making use of our choice of target coding scheme for the  $t_n$ , we obtain an expression for the bias in the form

$$w_0 = -\mathbf{w}^T \mathbf{m} \quad (4.34)$$

where we have used

$$\sum_{n=1}^N t_n = N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2} = 0 \quad (4.35)$$

and where  $\mathbf{m}$  is the mean of the total data set and is given by

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2). \quad (4.36)$$

*Exercise 4.6*

After some straightforward algebra, and again making use of the choice of  $t_n$ , the second equation (4.33) becomes

$$\left( \mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_2 - \mathbf{m}_1) \quad (4.37)$$

where  $\mathbf{S}_W$  is defined by (4.28),  $\mathbf{S}_B$  is defined by (4.27), and we have substituted for the bias using (4.34). Using (4.27), we note that  $\mathbf{S}_B \mathbf{w}$  is always in the direction of  $(\mathbf{m}_2 - \mathbf{m}_1)$ . Thus we can write

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \quad (4.38)$$

where we have ignored irrelevant scale factors. Thus the weight vector coincides with that found from the Fisher criterion. In addition, we have also found an expression for the bias value  $w_0$  given by (4.34). This tells us that a new vector  $\mathbf{x}$  should be classified as belonging to class  $\mathcal{C}_1$  if  $y(\mathbf{x}) = \mathbf{w}^T (\mathbf{x} - \mathbf{m}) > 0$  and class  $\mathcal{C}_2$  otherwise.

### 4.1.6 Fisher's discriminant for multiple classes

$D > K$   
但  $D > K$ ,  $K > D'$

We now consider the generalization of the Fisher discriminant to  $K > 2$  classes, and we shall assume that the dimensionality  $D$  of the input space is greater than the number  $K$  of classes. Next, we introduce  $D' > 1$  linear ‘features’  $y_k = \mathbf{w}_k^T \mathbf{x}$ , where  $k = 1, \dots, D'$ . These feature values can conveniently be grouped together to form a vector  $\mathbf{y}$ . Similarly, the weight vectors  $\{\mathbf{w}_k\}$  can be considered to be the columns of a matrix  $\mathbf{W}$ , so that

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}. \quad (4.39)$$

Note that again we are not including any bias parameters in the definition of  $\mathbf{y}$ . The generalization of the within-class covariance matrix to the case of  $K$  classes follows from (4.28) to give

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k \quad (4.40)$$

where

$$\mathbf{S}_k = \sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T \quad (4.41)$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n \quad (4.42)$$

and  $N_k$  is the number of patterns in class  $\mathcal{C}_k$ . In order to find a generalization of the between-class covariance matrix, we follow Duda and Hart (1973) and consider first the total covariance matrix

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T \quad (4.43)$$

where  $\mathbf{m}$  is the mean of the total data set

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k \quad (4.44)$$

and  $N = \sum_k N_k$  is the total number of data points. The total covariance matrix can be decomposed into the sum of the within-class covariance matrix, given by (4.40) and (4.41), plus an additional matrix  $\mathbf{S}_B$ , which we identify as a measure of the between-class covariance

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B \quad (4.45)$$

where

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T. \quad (4.46)$$

These covariance matrices have been defined in the original  $\mathbf{x}$ -space. We can now define similar matrices in the projected  $D'$ -dimensional  $\mathbf{y}$ -space

$$\textcolor{red}{S_W} \quad \mathbf{s}_{\mathbf{W}} = \sum_{k=1}^K \sum_{n \in C_k} (\mathbf{y}_n - \boldsymbol{\mu}_k)(\mathbf{y}_n - \boldsymbol{\mu}_k)^T \quad (4.47)$$

*大寫 "S" !*

and

$$\textcolor{red}{S_B} \quad \mathbf{s}_{\mathbf{B}} = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T \quad (4.48)$$

where

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}_n, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^K N_k \boldsymbol{\mu}_k. \quad (4.49)$$

Again we wish to construct a scalar that is large when the between-class covariance is large and when the within-class covariance is small. There are now many possible choices of criterion (Fukunaga, 1990). One example is given by

$$J(\mathbf{W}) = \text{Tr} \left\{ \frac{\mathbf{s}_B^T \mathbf{s}_B}{\mathbf{s}_W^T \mathbf{s}_W} \right\}. \quad (4.50)$$

This criterion can then be rewritten as an explicit function of the projection matrix  $\mathbf{W}$  in the form

$$J(\mathbf{W}) = \text{Tr} \left\{ (\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} (\mathbf{W}^T \mathbf{S}_B \mathbf{W}) \right\}. \quad (4.51)$$

Maximization of such criteria is straightforward, though somewhat involved, and is discussed at length in Fukunaga (1990). The weight values are determined by those eigenvectors of  $\mathbf{S}_W^{-1} \mathbf{S}_B$  that correspond to the  $D'$  largest eigenvalues.

There is one important result that is common to all such criteria, which is worth emphasizing. We first note from (4.46) that  $\mathbf{S}_B$  is composed of the sum of  $K$  matrices, each of which is an outer product of two vectors and therefore of rank 1. In addition, only  $(K - 1)$  of these matrices are independent as a result of the constraint (4.44). Thus,  $\mathbf{S}_B$  has rank at most equal to  $(K - 1)$  and so there are at most  $(K - 1)$  nonzero eigenvalues. This shows that the projection onto the  $(K - 1)$ -dimensional subspace spanned by the eigenvectors of  $\mathbf{S}_B$  does not alter the value of  $J(\mathbf{W})$ , and so we are therefore unable to find more than  $(K - 1)$  linear ‘features’ by this means (Fukunaga, 1990).

*K-1 > D'*

#### 4.1.7 The perceptron algorithm

Another example of a linear discriminant model is the perceptron of Rosenblatt (1962), which occupies an important place in the history of pattern recognition algorithms. It corresponds to a two-class model in which the input vector  $\mathbf{x}$  is first transformed using a fixed nonlinear transformation to give a feature vector  $\phi(\mathbf{x})$ , and this is then used to construct a generalized linear model of the form

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x})) \quad (4.52)$$

where the nonlinear activation function  $f(\cdot)$  is given by a step function of the form

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0. \end{cases} \quad (4.53)$$

The vector  $\phi(\mathbf{x})$  will typically include a bias component  $\phi_0(\mathbf{x}) = 1$ . In earlier discussions of two-class classification problems, we have focussed on a target coding scheme in which  $t \in \{0, 1\}$ , which is appropriate in the context of probabilistic models. For the perceptron, however, it is more convenient to use target values  $t = +1$  for class  $C_1$  and  $t = -1$  for class  $C_2$ , which matches the choice of activation function.

The algorithm used to determine the parameters  $\mathbf{w}$  of the perceptron can most easily be motivated by error function minimization. A natural choice of error function would be the total number of misclassified patterns. However, this does not lead to a simple learning algorithm because the error is a piecewise constant function of  $\mathbf{w}$ , with discontinuities wherever a change in  $\mathbf{w}$  causes the decision boundary to move across one of the data points. Methods based on changing  $\mathbf{w}$  using the gradient of the error function cannot then be applied, because the gradient is zero almost everywhere.

We therefore consider an alternative error function known as the *perceptron criterion*. To derive this, we note that we are seeking a weight vector  $\mathbf{w}$  such that patterns  $\mathbf{x}_n$  in class  $C_1$  will have  $\mathbf{w}^T \phi(\mathbf{x}_n) > 0$ , whereas patterns  $\mathbf{x}_n$  in class  $C_2$  have  $\mathbf{w}^T \phi(\mathbf{x}_n) < 0$ . Using the  $t \in \{-1, +1\}$  target coding scheme it follows that we would like all patterns to satisfy  $\mathbf{w}^T \phi(\mathbf{x}_n) t_n > 0$ . The perceptron criterion associates zero error with any pattern that is correctly classified, whereas for a misclassified pattern  $\mathbf{x}_n$  it tries to minimize the quantity  $-\mathbf{w}^T \phi(\mathbf{x}_n) t_n$ . The perceptron criterion is therefore given by

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n \quad (4.54)$$



**Frank Rosenblatt**  
1928–~~1969~~ 1971

Rosenblatt's perceptron played an important role in the history of machine learning. Initially, Rosenblatt simulated the perceptron on an IBM 704 computer at Cornell in 1957, but by the early 1960s he had built special-purpose hardware that provided a direct, parallel implementation of perceptron learning. Many of his ideas were encapsulated in "Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms" published in 1962. Rosenblatt's work was criticized by Marvin Minsky, whose objections were published in the book "Perceptrons", co-authored with

Seymour Papert. This book was widely misinterpreted at the time as showing that neural networks were fatally flawed and could only learn solutions for linearly separable problems. In fact, it only proved such limitations in the case of single-layer networks such as the perceptron and merely conjectured (incorrectly) that they applied to more general network models. Unfortunately, however, this book contributed to the substantial decline in research funding for neural computing, a situation that was not reversed until the mid-1980s. Today, there are many hundreds, if not thousands, of applications of neural networks in widespread use, with examples in areas such as handwriting recognition and information retrieval being used routinely by millions of people.

$$\phi_n = \phi(x_n) \text{ and}$$

where  $\mathcal{M}$  denotes the set of all misclassified patterns. The contribution to the error associated with a particular misclassified pattern is a linear function of  $\mathbf{w}$  in regions of  $\mathbf{w}$  space where the pattern is misclassified and zero in regions where it is correctly classified. The total error function is therefore piecewise linear.

### Section 3.1.3

We now apply the stochastic gradient descent algorithm to this error function. The change in the weight vector  $\mathbf{w}$  is then given by

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n \quad (4.55)$$

where  $\eta$  is the learning rate parameter and  $\tau$  is an integer that indexes the steps of the algorithm. Because the perceptron function  $y(\mathbf{x}, \mathbf{w})$  is unchanged if we multiply  $\mathbf{w}$  by a constant, we can set the learning rate parameter  $\eta$  equal to 1 without loss of generality. Note that, as the weight vector evolves during training, the set of patterns that are misclassified will change.

The perceptron learning algorithm has a simple interpretation, as follows. We cycle through the training patterns in turn, and for each pattern  $\mathbf{x}_n$  we evaluate the perceptron function (4.52). If the pattern is correctly classified, then the weight vector remains unchanged, whereas if it is incorrectly classified, then for class  $C_1$  we add the vector  $\phi(\mathbf{x}_n)$  onto the current estimate of weight vector  $\mathbf{w}$  while for class  $C_2$  we subtract the vector  $\phi(\mathbf{x}_n)$  from  $\mathbf{w}$ . The perceptron learning algorithm is illustrated in Figure 4.7.

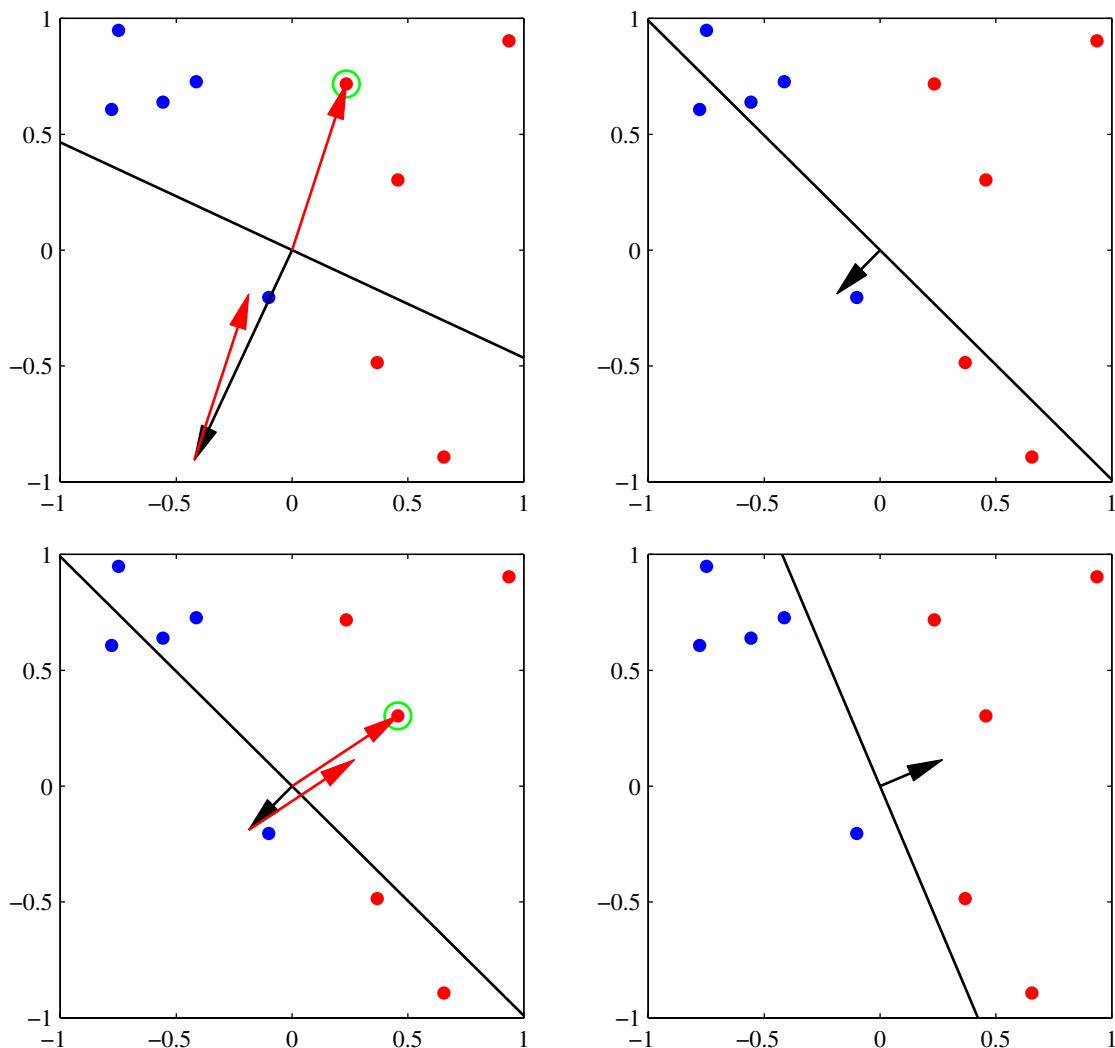
If we consider the effect of a single update in the perceptron learning algorithm, we see that the contribution to the error from a misclassified pattern will be reduced because from (4.55) we have

$$-\mathbf{w}^{(\tau+1)\top} \phi_n t_n = -\mathbf{w}^{(\tau)\top} \phi_n t_n - (\phi_n t_n)^\top \phi_n t_n < -\mathbf{w}^{(\tau)\top} \phi_n t_n \quad (4.56)$$

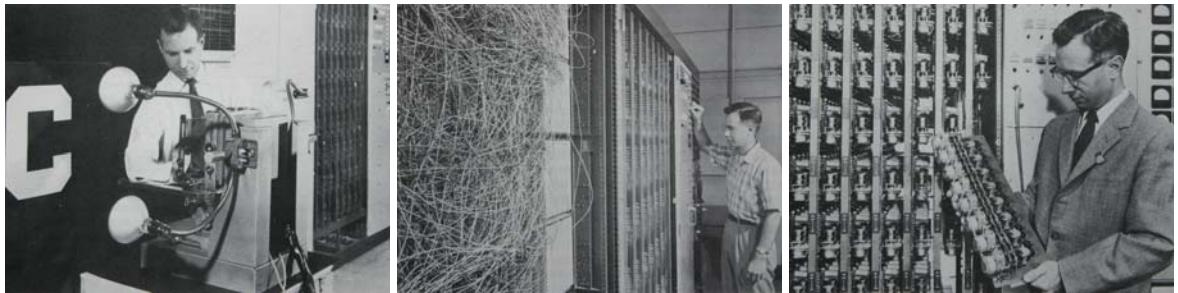
where we have set  $\eta = 1$ , and made use of  $\|\phi_n t_n\|^2 > 0$ . Of course, this does not imply that the contribution to the error function from the other misclassified patterns will have been reduced. Furthermore, the change in weight vector may have caused some previously correctly classified patterns to become misclassified. Thus the perceptron learning rule is not guaranteed to reduce the total error function at each stage.

However, the *perceptron convergence theorem* states that if there exists an exact solution (in other words, if the training data set is linearly separable), then the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps. Proofs of this theorem can be found for example in Rosenblatt (1962), Block (1962), Nilsson (1965), Minsky and Papert (1969), Hertz *et al.* (1991), and Bishop (1995a). Note, however, that the number of steps required to achieve convergence could still be substantial, and in practice, until convergence is achieved, we will not be able to distinguish between a nonseparable problem and one that is simply slow to converge.

Even when the data set is linearly separable, there may be many solutions, and which one is found will depend on the initialization of the parameters and on the order of presentation of the data points. Furthermore, for data sets that are not linearly separable, the perceptron learning algorithm will never converge.



**Figure 4.7** Illustration of the convergence of the perceptron learning algorithm, showing data points from two classes (red and blue) in a two-dimensional feature space ( $\phi_1, \phi_2$ ). The top left plot shows the initial parameter vector  $w$  shown as a black arrow together with the corresponding decision boundary (black line), in which the arrow points towards the decision region which classified as belonging to the red class. The data point circled in green is misclassified and so its feature vector is added to the current weight vector, giving the new decision boundary shown in the top right plot. The bottom left plot shows the next misclassified point to be considered, indicated by the green circle, and its feature vector is again added to the weight vector giving the decision boundary shown in the bottom right plot for which all data points are correctly classified.



**Figure 4.8** Illustration of the Mark 1 perceptron hardware. The photograph on the left shows how the inputs were obtained using a simple camera system in which an input scene, in this case a printed character, was illuminated by powerful lights, and an image focussed onto a  $20 \times 20$  array of cadmium sulphide photocells, giving a primitive 400 pixel image. The perceptron also had a patch board, shown in the middle photograph, which allowed different configurations of input features to be tried. Often these were wired up at random to demonstrate the ability of the perceptron to learn without the need for precise wiring, in contrast to a modern digital computer. The photograph on the right shows one of the racks of adaptive weights. Each weight was implemented using a rotary variable resistor, also called a potentiometer, driven by an electric motor thereby allowing the value of the weight to be adjusted automatically by the learning algorithm.

Aside from difficulties with the learning algorithm, the perceptron does not provide probabilistic outputs, nor does it generalize readily to  $K > 2$  classes. The most important limitation, however, arises from the fact that (in common with all of the models discussed in this chapter and the previous one) it is based on linear combinations of fixed basis functions. More detailed discussions of the limitations of perceptrons can be found in Minsky and Papert (1969) and Bishop (1995a).

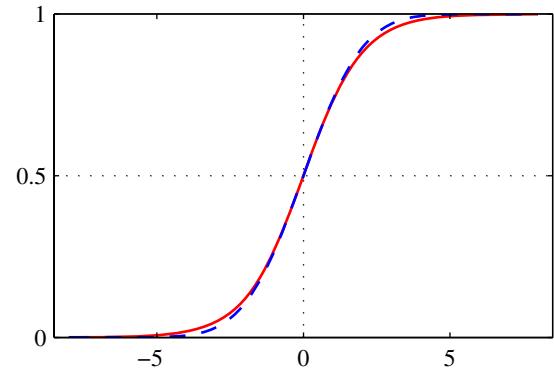
Analogue hardware implementations of the perceptron were built by Rosenblatt, based on motor-driven variable resistors to implement the adaptive parameters  $w_j$ . These are illustrated in Figure 4.8. The inputs were obtained from a simple camera system based on an array of photo-sensors, while the basis functions  $\phi$  could be chosen in a variety of ways, for example based on simple fixed functions of randomly chosen subsets of pixels from the input image. Typical applications involved learning to discriminate simple shapes or characters.

At the same time that the perceptron was being developed, a closely related system called the *adaline*, which is short for ‘adaptive linear element’, was being explored by Widrow and co-workers. The functional form of the model was the same as for the perceptron, but a different approach to training was adopted (Widrow and Hoff, 1960; Widrow and Lehr, 1990).

## 4.2. Probabilistic Generative Models

We turn next to a probabilistic view of classification and show how models with linear decision boundaries arise from simple assumptions about the distribution of the data. In Section 1.5.4, we discussed the distinction between the discriminative and the generative approaches to classification. Here we shall adopt a generative

**Figure 4.9** Plot of the logistic sigmoid function  $\sigma(a)$  defined by (4.59), shown in red, together with the scaled probit function  $\Phi(\lambda a)$ , for  $\lambda^2 = \pi/8$ , shown in dashed blue, where  $\Phi(a)$  is defined by (4.114). The scaling factor  $\pi/8$  is chosen so that the derivatives of the two curves are equal for  $a = 0$ .



approach in which we model the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ , as well as the class priors  $p(\mathcal{C}_k)$ , and then use these to compute posterior probabilities  $p(\mathcal{C}_k|\mathbf{x})$  through Bayes' theorem.

Consider first of all the case of two classes. The posterior probability for class  $\mathcal{C}_1$  can be written as

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned} \quad (4.57)$$

where we have defined

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \quad (4.58)$$

and  $\sigma(a)$  is the *logistic sigmoid* function defined by

$$\sigma(a) = \frac{1}{1 + \exp(-a)} \quad (4.59)$$

which is plotted in Figure 4.9. The term ‘sigmoid’ means S-shaped. This type of function is sometimes also called a ‘squashing function’ because it maps the whole real axis into a finite interval. The logistic sigmoid has been encountered already in earlier chapters and plays an important role in many classification algorithms. It satisfies the following symmetry property

$$\sigma(-a) = 1 - \sigma(a) \quad (4.60)$$

as is easily verified. The inverse of the logistic sigmoid is given by

$$a = \ln \left( \frac{\sigma}{1 - \sigma} \right) \quad (4.61)$$

and is known as the *logit* function. It represents the log of the ratio of probabilities  $[p(\mathcal{C}_1|\mathbf{x})/p(\mathcal{C}_2|\mathbf{x})]$  for the two classes, also known as the *log odds*.

Note that in (4.57) we have simply rewritten the posterior probabilities in an equivalent form, and so the appearance of the logistic sigmoid may seem rather vacuous. However, it will have significance provided  $a(\mathbf{x})$  takes a simple functional form. We shall shortly consider situations in which  $a(\mathbf{x})$  is a linear function of  $\mathbf{x}$ , in which case the posterior probability is governed by a generalized linear model.

For the case of  $K > 2$  classes, we have

$$\begin{aligned} p(\mathcal{C}_k|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned} \quad (4.62)$$

which is known as the *normalized exponential* and can be regarded as a multiclass generalization of the logistic sigmoid. Here the quantities  $a_k$  are defined by

$$a_k = \ln(p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)) \quad (4.63)$$

The normalized exponential is also known as the *softmax function*, as it represents a smoothed version of the ‘max’ function because, if  $a_k \gg a_j$  for all  $j \neq k$ , then  $p(\mathcal{C}_k|\mathbf{x}) \simeq 1$ , and  $p(\mathcal{C}_j|\mathbf{x}) \simeq 0$ .

We now investigate the consequences of choosing specific forms for the class-conditional densities, looking first at continuous input variables  $\mathbf{x}$  and then discussing briefly the case of discrete inputs.

### 4.2.1 Continuous inputs 从高斯分布到逻辑斯谛回归

Let us assume that the class-conditional densities are Gaussian and then explore the resulting form for the posterior probabilities. To start with, we shall assume that all classes share the same covariance matrix. Thus the density for class  $\mathcal{C}_k$  is given by

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}. \quad (4.64)$$

**二分类情况** Consider first the case of two classes. From (4.57) and (4.58), we have

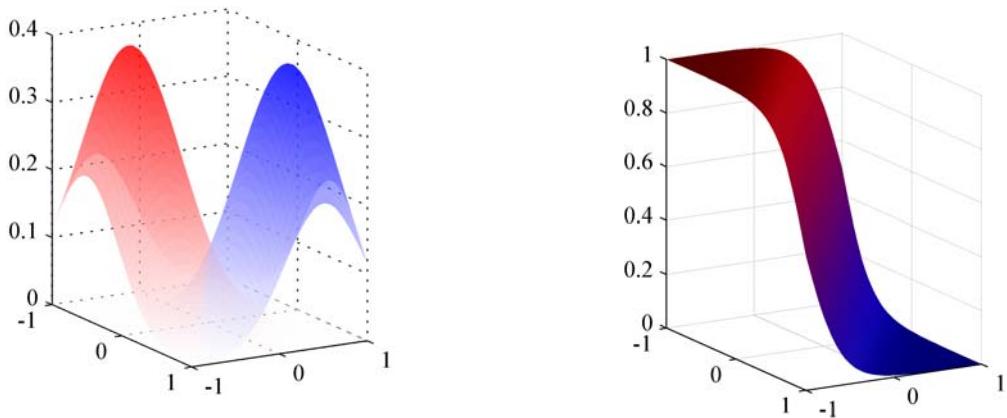
$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \quad (4.65)$$

where we have defined

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad (4.66)$$

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}. \quad (4.67)$$

We see that the quadratic terms in  $\mathbf{x}$  from the exponents of the Gaussian densities have cancelled (due to the assumption of common covariance matrices) leading to a linear function of  $\mathbf{x}$  in the argument of the logistic sigmoid. This result is illustrated for the case of a two-dimensional input space  $\mathbf{x}$  in Figure 4.10. The resulting



**Figure 4.10** The left-hand plot shows the class-conditional densities for two classes, denoted red and blue. On the right is the corresponding posterior probability  $p(\mathcal{C}_1|\mathbf{x})$ , which is given by a logistic sigmoid of a linear function of  $\mathbf{x}$ . The surface in the right-hand plot is coloured using a proportion of red ink given by  $p(\mathcal{C}_1|\mathbf{x})$  and a proportion of blue ink given by  $p(\mathcal{C}_2|\mathbf{x}) = 1 - p(\mathcal{C}_1|\mathbf{x})$ .

decision boundaries correspond to surfaces along which the posterior probabilities  $p(\mathcal{C}_k|\mathbf{x})$  are constant and so will be given by linear functions of  $\mathbf{x}$ , and therefore the decision boundaries are linear in input space. The prior probabilities  $p(\mathcal{C}_k)$  enter only through the bias parameter  $w_0$  so that changes in the priors have the effect of making parallel shifts of the decision boundary and more generally of the parallel contours of constant posterior probability.]

For the general case of  $K$  classes we have, from (4.62) and (4.63),

$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (4.68)$$

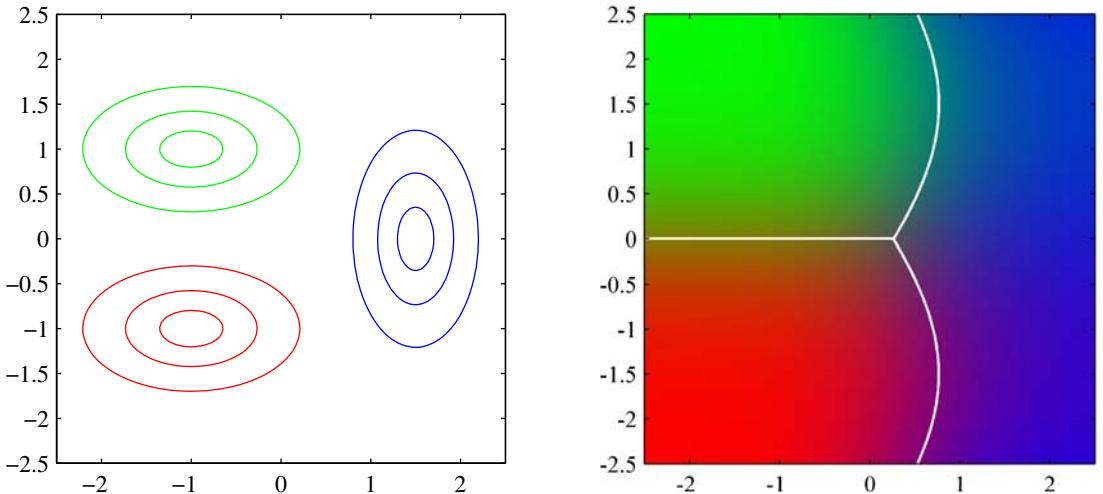
where we have defined

$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k \quad (4.69)$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k). \quad (4.70)$$

We see that the  $a_k(\mathbf{x})$  are again linear functions of  $\mathbf{x}$  as a consequence of the cancellation of the quadratic terms due to the shared covariances.] The resulting decision boundaries, corresponding to the minimum misclassification rate, will occur when two of the posterior probabilities (the two largest) are equal, and so will be defined by linear functions of  $\mathbf{x}$ , and so again we have a generalized linear model.]

If we relax the assumption of a shared covariance matrix and allow each class-conditional density  $p(\mathbf{x}|\mathcal{C}_k)$  to have its own covariance matrix  $\Sigma_k$ , then the earlier cancellations will no longer occur, and we will obtain quadratic functions of  $\mathbf{x}$ , giving rise to a *quadratic discriminant*. The linear and quadratic decision boundaries are illustrated in Figure 4.11.



**Figure 4.11** The left-hand plot shows the class-conditional densities for three classes each having a Gaussian distribution, coloured red, green, and blue, in which the red and green classes have the same covariance matrix. The right-hand plot shows the corresponding posterior probabilities, in which the RGB colour vector represents the posterior probabilities for the respective three classes. The decision boundaries are also shown. Notice that the boundary between the red and green classes, which have the same covariance matrix, is linear, whereas those between the other pairs of classes are quadratic.

### 4.2.2 Maximum likelihood solution

Once we have specified a parametric functional form for the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ , we can then determine the values of the parameters, together with the prior class probabilities  $p(\mathcal{C}_k)$ , using maximum likelihood. This requires a data set comprising observations of  $\mathbf{x}$  along with their corresponding class labels.

Consider first the case of two classes, each having a Gaussian class-conditional density with a shared covariance matrix, and suppose we have a data set  $\{\mathbf{x}_n, t_n\}$  where  $n = 1, \dots, N$ . Here  $t_n = 1$  denotes class  $\mathcal{C}_1$  and  $t_n = 0$  denotes class  $\mathcal{C}_2$ . We denote the prior class probability  $p(\mathcal{C}_1) = \pi$ , so that  $p(\mathcal{C}_2) = 1 - \pi$ . For a data point  $\mathbf{x}_n$  from class  $\mathcal{C}_1$ , we have  $t_n = 1$  and hence

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1)p(\mathbf{x}_n|\mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}).$$

Similarly for class  $\mathcal{C}_2$ , we have  $t_n = 0$  and hence

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2)p(\mathbf{x}_n|\mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}).$$

Thus the likelihood function is given by

$$\text{lik}(\mathbf{t}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n} \quad (4.71)$$

where  $\mathbf{t} = (t_1, \dots, t_N)^T$ . As usual, it is convenient to maximize the log of the likelihood function. Consider first the maximization with respect to  $\pi$ . The terms in

the log likelihood function that depend on  $\pi$  are

$$\sum_{n=1}^N \{t_n \ln \pi + (1 - t_n) \ln(1 - \pi)\}. \quad (4.72)$$

Setting the derivative with respect to  $\pi$  equal to zero and rearranging, we obtain

$$\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} \quad (4.73)$$

where  $N_1$  denotes the total number of data points in class  $\mathcal{C}_1$ , and  $N_2$  denotes the total number of data points in class  $\mathcal{C}_2$ . Thus the maximum likelihood estimate for  $\pi$  is simply the fraction of points in class  $\mathcal{C}_1$  as expected. This result is easily generalized to the multiclass case where again the maximum likelihood estimate of the prior probability associated with class  $\mathcal{C}_k$  is given by the fraction of the training set points assigned to that class.

### Exercise 4.9

Now consider the maximization with respect to  $\boldsymbol{\mu}_1$ . Again we can pick out of the log likelihood function those terms that depend on  $\boldsymbol{\mu}_1$  giving

$$\sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) + \text{const.} \quad (4.74)$$

Setting the derivative with respect to  $\boldsymbol{\mu}_1$  to zero and rearranging, we obtain

$$\boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n \quad (4.75)$$

which is simply the mean of all the input vectors  $\mathbf{x}_n$  assigned to class  $\mathcal{C}_1$ . By a similar argument, the corresponding result for  $\boldsymbol{\mu}_2$  is given by

$$\boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n \quad (4.76)$$

which again is the mean of all the input vectors  $\mathbf{x}_n$  assigned to class  $\mathcal{C}_2$ .

Finally, consider the maximum likelihood solution for the shared covariance matrix  $\boldsymbol{\Sigma}$ . Picking out the terms in the log likelihood function that depend on  $\boldsymbol{\Sigma}$ , we have

$$\begin{aligned} & -\frac{1}{2} \sum_{n=1}^N t_n \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \\ & - \frac{1}{2} \sum_{n=1}^N (1 - t_n) \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (1 - t_n) (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \\ & = -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{N}{2} \text{Tr} \{ \boldsymbol{\Sigma}^{-1} \mathbf{S} \} \end{aligned} \quad (4.77)$$

where we have defined

$$\mathbf{S} = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2 \quad (4.78)$$

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \quad (4.79)$$

$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T. \quad (4.80)$$

Using the standard result for the maximum likelihood solution for a Gaussian distribution, we see that  $\Sigma = \mathbf{S}$ , which represents a weighted average of the covariance matrices associated with each of the two classes separately.

This result is easily extended to the  $K$  class problem to obtain the corresponding maximum likelihood solutions for the parameters in which each class-conditional density is Gaussian with a shared covariance matrix. Note that the approach of fitting Gaussian distributions to the classes is not robust to outliers, because the maximum likelihood estimation of a Gaussian is not robust.

### Exercise 4.10

#### Section 2.3.7

### 朴素贝叶斯假设

#### Section 8.2.2

### Exercise 4.11

总结前文并给出一般结论：  
由指数族分布推得  
逻辑斯谛回归。

### 4.2.3 Discrete features 从离散变量的朴素贝叶斯假设到逻辑斯谛回归

Let us now consider the case of discrete feature values  $x_i$ . For simplicity, we begin by looking at binary feature values  $x_i \in \{0, 1\}$  and discuss the extension to more general discrete features shortly. If there are  $D$  inputs, then a general distribution would correspond to a table of  $2^D$  numbers for each class, containing  $2^D - 1$  independent variables (due to the summation constraint). Because this grows exponentially with the number of features, we might seek a more restricted representation. Here we will make the *naive Bayes* assumption in which the feature values are treated as independent, conditioned on the class  $\mathcal{C}_k$ . Thus we have class-conditional distributions of the form

$$p(\mathbf{x}|\mathcal{C}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \quad (4.81)$$

which contain  $D$  independent parameters for each class. Substituting into (4.63) then gives

$$a_k(\mathbf{x}) = \sum_{i=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln(1 - \mu_{ki})\} + \ln p(\mathcal{C}_k) \quad (4.82)$$

which again are linear functions of the input values  $x_i$ . For the case of  $K = 2$  classes, we can alternatively consider the logistic sigmoid formulation given by (4.57). Analogous results are obtained for discrete variables each of which can take  $M > 2$  states.

### 4.2.4 Exponential family

As we have seen, for both Gaussian distributed and discrete inputs, the posterior class probabilities are given by generalized linear models with logistic sigmoid ( $K =$

2 classes) or softmax ( $K \geq 2$  classes) activation functions. These are particular cases of a more general result obtained by assuming that the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$  are members of the exponential family of distributions. J

证明: Using the form (2.194) for members of the exponential family, we see that the distribution of  $\mathbf{x}$  can be written in the form

$$p(\mathbf{x}|\boldsymbol{\lambda}_k) = h(\mathbf{x})g(\boldsymbol{\lambda}_k) \exp \left\{ \boldsymbol{\lambda}_k^T \mathbf{u}(\mathbf{x}) \right\}. \quad (4.83)$$

We now restrict attention to the subclass of such distributions for which  $\mathbf{u}(\mathbf{x}) = \mathbf{x}$ . Then we make use of (2.236) to introduce a scaling parameter  $s$ , so that we obtain the restricted set of exponential family class-conditional densities of the form

$$p(\mathbf{x}|\boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s}\mathbf{x}\right) g(\boldsymbol{\lambda}_k) \exp \left\{ \frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x} \right\}. \quad (4.84)$$

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  $s$ .

For the two-class problem, we substitute this expression for the class-conditional densities into (4.58) and we see that the posterior class probability is again given by a logistic sigmoid acting on a linear function  $a(\mathbf{x})$  which is given by

$$a(\mathbf{x}) = \frac{1}{s} (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_1) - \ln g(\boldsymbol{\lambda}_2) + \ln p(\mathcal{C}_1) - \ln p(\mathcal{C}_2). \quad (4.85)$$

Similarly, for the  $K$ -class problem, we substitute the class-conditional density expression into (4.63) to give

$$a_k(\mathbf{x}) = \frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_k) + \ln p(\mathcal{C}_k) \quad (4.86)$$

and so again is a linear function of  $\mathbf{x}$ . J

### 4.3. Probabilistic Discriminative Models

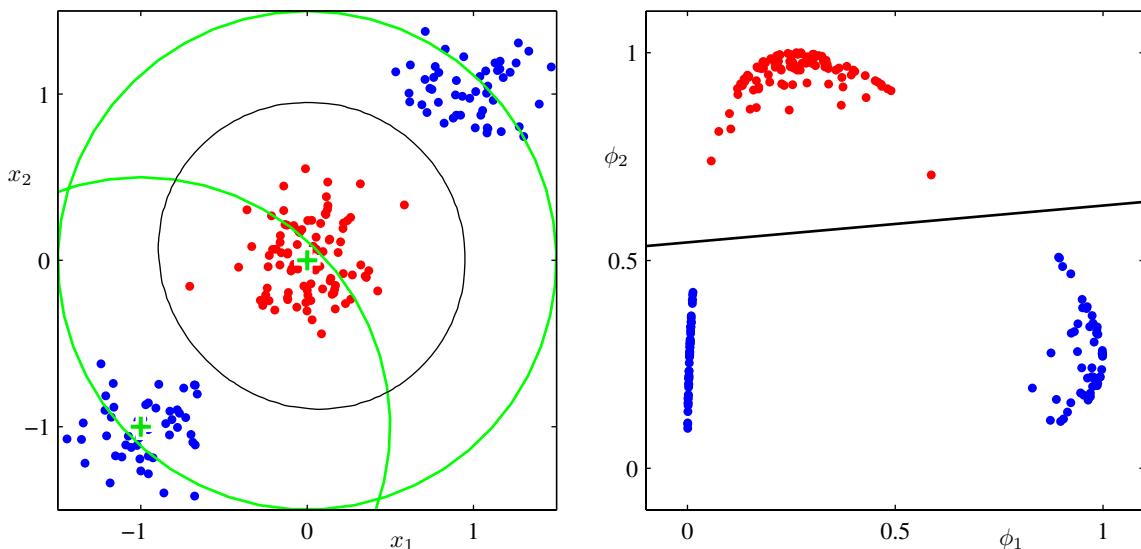
4.3节总结

For the two-class classification problem, we have seen that the posterior probability of class  $\mathcal{C}_1$  can be written as a logistic sigmoid acting on a linear function of  $\mathbf{x}$ , for a wide choice of class-conditional distributions  $p(\mathbf{x}|\mathcal{C}_k)$ . Similarly, for the multiclass case, the posterior probability of class  $\mathcal{C}_k$  is given by a softmax transformation of a linear function of  $\mathbf{x}$ . For specific choices of the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ , we have used maximum likelihood to determine the parameters of the densities as well as the class priors  $p(\mathcal{C}_k)$  and then used Bayes' theorem to find the posterior class probabilities.

进一步引出4.3节的思路 However, an alternative approach is to use the functional form of the generalized linear model explicitly and to determine its parameters directly by using maximum likelihood. We shall see that there is an efficient algorithm finding such solutions known as *iterative reweighted least squares*, or *IRLS*.

The indirect approach to finding the parameters of a generalized linear model, by fitting class-conditional densities and class priors separately and then applying

阐述生成模型及判别模型的含义，并比较优缺点。



**Figure 4.12** Illustration of the role of nonlinear basis functions in linear classification models. The left plot shows the original input space  $(x_1, x_2)$  together with data points from two classes labelled red and blue. Two ‘Gaussian’ basis functions  $\phi_1(\mathbf{x})$  and  $\phi_2(\mathbf{x})$  are defined in this space with centres shown by the green crosses and with contours shown by the green circles. The right-hand plot shows the corresponding feature space  $(\phi_1, \phi_2)$  together with the linear decision boundary obtained given by a logistic regression model of the form discussed in Section 4.3.2. This corresponds to a nonlinear decision boundary in the original input space, shown by the black curve in the left-hand plot.

Bayes’ theorem, represents an example of *generative modelling*, because we could take such a model and generate synthetic data by drawing values of  $\mathbf{x}$  from the marginal distribution  $p(\mathbf{x})$ . In the direct approach, we are maximizing a likelihood function defined through the conditional distribution  $p(C_k|\mathbf{x})$ , which represents a form of *discriminative training*. One advantage of the discriminative approach is that there will typically be fewer adaptive parameters to be determined, as we shall see shortly. It may also lead to improved predictive performance, particularly when the class-conditional density assumptions give a poor approximation to the true distributions.

### 4.3.1 Fixed basis functions

So far in this chapter, we have considered classification models that work directly with the original input vector  $\mathbf{x}$ . However, all of the algorithms are equally applicable if we first make a fixed nonlinear transformation of the inputs using a vector of basis functions  $\phi(\mathbf{x})$ . The resulting decision boundaries will be linear in the feature space  $\phi$ , and these correspond to nonlinear decision boundaries in the original  $\mathbf{x}$  space, as illustrated in Figure 4.12. Classes that are linearly separable in the feature space  $\phi(\mathbf{x})$  need not be linearly separable in the original observation space  $\mathbf{x}$ . Note that as in our discussion of linear models for regression, one of the

basis functions is typically set to a constant, say  $\phi_0(\mathbf{x}) = 1$ , so that the corresponding parameter  $w_0$  plays the role of a bias. For the remainder of this chapter, we shall include a fixed basis function transformation  $\phi(\mathbf{x})$ , as this will highlight some useful similarities to the regression models discussed in Chapter 3.

For many problems of practical interest, there is significant overlap between the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ . This corresponds to posterior probabilities  $p(\mathcal{C}_k|\mathbf{x})$ , which, for at least some values of  $\mathbf{x}$ , are not 0 or 1. In such cases, the optimal solution is obtained by modelling the posterior probabilities accurately and then applying standard decision theory, as discussed in Chapter 1. Note that nonlinear transformations  $\phi(\mathbf{x})$  cannot remove such class overlap. Indeed, they can increase the level of overlap, or create overlap where none existed in the original observation space. However, suitable choices of nonlinearity can make the process of modelling the posterior probabilities easier.

### Section 3.6

Such fixed basis function models have important limitations, and these will be resolved in later chapters by allowing the basis functions themselves to adapt to the data. Notwithstanding these limitations, models with fixed nonlinear basis functions play an important role in applications, and a discussion of such models will introduce many of the key concepts needed for an understanding of their more complex counterparts.

#### 4.3.2 Logistic regression

We begin our treatment of generalized linear models by considering the problem of two-class classification. In our discussion of generative approaches in Section 4.2, we saw that under rather general assumptions, the posterior probability of class  $\mathcal{C}_1$  can be written as a logistic sigmoid acting on a linear function of the feature vector  $\phi$  so that

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (4.87)$$

with  $p(\mathcal{C}_2|\phi) = 1 - p(\mathcal{C}_1|\phi)$ . Here  $\sigma(\cdot)$  is the *logistic sigmoid* function defined by (4.59). In the terminology of statistics, this model is known as *logistic regression*, although it should be emphasized that this is a model for classification rather than regression.

*手稿文*

For an  $M$ -dimensional feature space  $\phi$ , this model has  $M$  adjustable parameters. By contrast, if we had fitted Gaussian class conditional densities using maximum likelihood, we would have used  $2M$  parameters for the means and  $M(M+1)/2$  parameters for the (shared) covariance matrix. Together with the class prior  $p(\mathcal{C}_1)$ , this gives a total of  $M(M+5)/2 + 1$  parameters, which grows quadratically with  $M$ , in contrast to the linear dependence on  $M$  of the number of parameters in logistic regression. For large values of  $M$ , there is a clear advantage in working with the logistic regression model directly.

We now use maximum likelihood to determine the parameters of the logistic regression model. To do this, we shall make use of the derivative of the logistic sigmoid function, which can conveniently be expressed in terms of the sigmoid function itself

$$\frac{d\sigma}{da} = \sigma(1 - \sigma). \quad (4.88)$$

### Exercise 4.12

For a data set  $\{\phi_n, t_n\}$ , where  $t_n \in \{0, 1\}$  and  $\phi_n = \phi(\mathbf{x}_n)$ , with  $n = 1, \dots, N$ , the likelihood function can be written

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad (4.89)$$

where  $\mathbf{t} = (t_1, \dots, t_N)^T$  and  $y_n = p(\mathcal{C}_1|\phi_n)$ . As usual, we can define an error function by taking the negative logarithm of the likelihood, which gives the *cross-entropy* error function in the form

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (4.90)$$

**Exercise 4.13** where  $y_n = \sigma(a_n)$  and  $a_n = \mathbf{w}^T \phi_n$ . Taking the gradient of the error function with respect to  $\mathbf{w}$ , we obtain

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n \quad (4.91)$$

where we have made use of (4.88). We see that the factor involving the derivative of the logistic sigmoid has cancelled, leading to a simplified form for the gradient of the log likelihood. In particular, the contribution to the gradient from data point  $n$  is given by the ‘error’  $y_n - t_n$  between the target value and the prediction of the model, times the basis function vector  $\phi_n$ . Furthermore, comparison with (3.13) shows that this takes precisely the same form as the gradient of the sum-of-squares error function for the linear regression model.

### Section 3.1.1

**Sequential learning** If desired, we could make use of the result (4.91) to give a sequential algorithm in which patterns are presented one at a time, in which each of the weight vectors is updated using (3.22) in which  $\nabla E_n$  is the  $n^{\text{th}}$  term in (4.91).

**EXAMPLE** **是否存在问题:** It is worth noting that maximum likelihood can exhibit severe over-fitting for data sets that are linearly separable. This arises because the maximum likelihood solution occurs when the hyperplane corresponding to  $\sigma = 0.5$ , equivalent to  $\mathbf{w}^T \phi = 0$ , separates the two classes and the magnitude of  $\mathbf{w}$  goes to infinity. In this case, the logistic sigmoid function becomes infinitely steep in feature space, corresponding to a Heaviside step function, so that every training point from each class  $k$  is assigned a posterior probability  $p(\mathcal{C}_k|\mathbf{x}) = 1$ . Furthermore, there is typically a continuum of such solutions because any separating hyperplane will give rise to the same posterior probabilities at the training data points, as will be seen later in Figure 10.13.

Maximum likelihood provides no way to favour one such solution over another, and which solution is found in practice will depend on the choice of optimization algorithm and on the parameter initialization. Note that the problem will arise even if the number of data points is large compared with the number of parameters in the model, so long as the training data set is linearly separable. The singularity can be avoided by inclusion of a prior and finding a MAP solution for  $\mathbf{w}$ , or equivalently by adding a regularization term to the error function.

### Exercise 4.14

### 4.3.3 Iterative reweighted least squares

In the case of the linear regression models discussed in Chapter 3, the maximum likelihood solution, on the assumption of a Gaussian noise model, leads to a closed-form solution. This was a consequence of the quadratic dependence of the log likelihood function on the parameter vector  $\mathbf{w}$ . For logistic regression, there is no longer a closed-form solution, due to the nonlinearity of the logistic sigmoid function. However, the departure from a quadratic form is not substantial. To be precise, the error function is ~~convex~~, as we shall see shortly, and hence has a unique minimum. Furthermore, the error function can be minimized by an efficient iterative technique based on the *Newton-Raphson* iterative optimization scheme, which uses a local quadratic approximation to the log likelihood function. The Newton-Raphson update, for minimizing a function  $E(\mathbf{w})$ , takes the form (Fletcher, 1987; Bishop and Nabney, 2008)

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w}). \quad (4.92)$$

where  $\mathbf{H}$  is the Hessian matrix whose elements comprise the second derivatives of  $E(\mathbf{w})$  with respect to the components of  $\mathbf{w}$ .

[ Let us first of all apply the Newton-Raphson method to the linear regression model (3.3) with the sum-of-squares error function (3.12). The gradient and Hessian of this error function are given by

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t} \quad (4.93)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi \quad (4.94)$$

#### Section 3.1.1

where  $\Phi$  is the  $N \times M$  design matrix, whose  $n^{\text{th}}$  row is given by  $\phi_n^T$ . The Newton-Raphson update then takes the form

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \Phi)^{-1} \{ \Phi^T \Phi \mathbf{w}^{(\text{old})} - \Phi^T \mathbf{t} \} \\ &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \end{aligned} \quad (4.95)$$

which we recognize as the standard least-squares solution. Note that the error function in this case is quadratic and hence the Newton-Raphson formula gives the exact solution in one step.]

Now let us apply the Newton-Raphson update to the cross-entropy error function (4.90) for the logistic regression model. From (4.91) we see that the gradient and Hessian of this error function are given by

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t}) \quad (4.96)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi \quad (4.97)$$

where we have made use of (4.88). Also, we have introduced the  $N \times N$  diagonal matrix  $\mathbf{R}$  with elements

$$R_{nn} = y_n(1 - y_n). \quad (4.98)$$

We see that the Hessian is no longer constant but depends on  $\mathbf{w}$  through the weighting matrix  $\mathbf{R}$ , corresponding to the fact that the error function is no longer quadratic. Using the property  $0 < y_n < 1$ , which follows from the form of the logistic sigmoid function, we see that  $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$  for an arbitrary vector  $\mathbf{u}$ , and so the Hessian matrix  $\mathbf{H}$  is positive definite. It follows that the error function is a **convex** function of  $\mathbf{w}$  and hence has a unique minimum. 这并不意味着  $\mathbf{w}$  的解只有一个

### Exercise 4.15

The Newton-Raphson update formula for the logistic regression model then becomes

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\text{old})} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z} \end{aligned} \quad (4.99)$$

where  $\mathbf{z}$  is an  $N$ -dimensional vector with elements

$$\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}). \quad (4.100)$$

*但这里并不介绍 Weighted least-squares 是什么*

We see that the update formula (4.99) takes the form of a set of normal equations for a weighted least-squares problem. Because the weighing matrix  $\mathbf{R}$  is not constant but depends on the parameter vector  $\mathbf{w}$ , we must apply the normal equations iteratively, each time using the new weight vector  $\mathbf{w}$  to compute a revised weighing matrix  $\mathbf{R}$ . For this reason, the algorithm is known as *iterative reweighted least squares*, or *IRLS* (Rubin, 1983). As in the weighted least-squares problem, the elements of the diagonal weighting matrix  $\mathbf{R}$  can be interpreted as variances because the mean and variance of  $t$  in the logistic regression model are given by

$$\mathbb{E}[t] = \sigma(\mathbf{x}) = y \quad (4.101)$$

$$\text{var}[t] = \mathbb{E}[t^2] - \mathbb{E}[t]^2 = \sigma(\mathbf{x}) - \sigma(\mathbf{x})^2 = y(1 - y) \quad (4.102)$$

where we have used the property  $t^2 = t$  for  $t \in \{0, 1\}$ . In fact, we can interpret IRLS as the solution to a linearized problem in the space of the variable  $a = \mathbf{w}^T \phi$ . The quantity  $z_n$ , which corresponds to the  $n^{\text{th}}$  element of  $\mathbf{z}$ , can then be given a simple interpretation as an effective target value in this space obtained by making a local linear approximation to the logistic sigmoid function around the current operating point  $\mathbf{w}^{(\text{old})}$

$$\begin{aligned} a_n(\mathbf{w}) &\simeq a_n(\mathbf{w}^{(\text{old})}) + \frac{da_n}{dy_n} \Big|_{\mathbf{w}^{(\text{old})}} (t_n - y_n) \\ &= \phi_n^T \mathbf{w}^{(\text{old})} - \frac{(y_n - t_n)}{y_n(1 - y_n)} = z_n. \end{aligned} \quad (4.103)$$

### 4.3.4 Multiclass logistic regression

*Section 4.2 4.2 产生模型* In our discussion of generative models for multiclass classification, we have seen that for a large class of distributions, the posterior probabilities are given by a softmax transformation of linear functions of the feature variables, so that

$$p(\mathcal{C}_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad (4.104)$$

where the ‘activations’  $a_k$  are given by

$$a_k = \mathbf{w}_k^T \phi. \quad (4.105)$$

There we used maximum likelihood to determine separately the class-conditional densities and the class priors and then found the corresponding posterior probabilities using Bayes’ theorem, thereby implicitly determining the parameters  $\{\mathbf{w}_k\}$ . Here we consider the use of maximum likelihood to determine the parameters  $\{\mathbf{w}_k\}$  of this model directly. To do this, we will require the derivatives of  $y_k$  with respect to all of the activations  $a_j$ . These are given by

$$\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j) \quad (4.106)$$

where  $I_{kj}$  are the elements of the identity matrix.

Next we write down the likelihood function. This is most easily done using the 1-of- $K$  coding scheme in which the target vector  $\mathbf{t}_n$  for a feature vector  $\phi_n$  belonging to class  $\mathcal{C}_k$  is a binary vector with all elements zero except for element  $k$ , which equals one. The likelihood function is then given by

$$p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(\mathcal{C}_k|\phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \quad (4.107)$$

where  $y_{nk} = y_k(\phi_n)$ , and  $\mathbf{T}$  is an  $N \times K$  matrix of target variables with elements  $t_{nk}$ . Taking the negative logarithm then gives

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \quad (4.108)$$

which is known as the *cross-entropy* error function for the multiclass classification problem.

We now take the gradient of the error function with respect to one of the parameter vectors  $\mathbf{w}_j$ . Making use of the result (4.106) for the derivatives of the softmax function, we obtain

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n \quad (4.109)$$

*Exercise 4.18*

where we have made use of  $\sum_k t_{nk} = 1$ . Once again, we see the same form arising for the gradient as was found for the sum-of-squares error function with the linear model and the cross-entropy error for the logistic regression model, namely the product of the error  $(y_{nj} - t_{nj})$  times the basis function  $\phi_n$ . Again, we could use this to formulate a sequential algorithm in which patterns are presented one at a time, in which each of the weight vectors is updated using (3.22).

We have seen that the derivative of the log likelihood function for a linear regression model with respect to the parameter vector  $w$  for a data point  $n$  took the form of the ‘error’  $y_n - t_n$  times the feature vector  $\phi_n$ . Similarly, for the combination of logistic sigmoid activation function and cross-entropy error function (4.90), and for the softmax activation function with the multiclass cross-entropy error function (4.108), we again obtain this same simple form. This is an example of a more general result, as we shall see in Section 4.3.6. 上述3种类型的梯度形式存在更一般化结论

To find a batch algorithm, we again appeal to the Newton-Raphson update to obtain the corresponding IRLS algorithm for the multiclass problem. This requires evaluation of the Hessian matrix that comprises blocks of size  $M \times M$  in which block  $j, k$  is given by

$$\nabla_{w_k} \nabla_{w_j} E(w_1, \dots, w_K) = \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T. \quad (4.110)$$

### Exercise 4.20

As with the two-class problem, the Hessian matrix for the multiclass logistic regression model is positive definite and so the error function again has a unique minimum. Practical details of IRLS for the multiclass case can be found in Bishop and Nabney (2008).

#### 4.3.5 Probit regression

对上文介绍的朴素  
判别模型进行总结

介绍一种有利于前  
文的新加权判  
别模型

[ We have seen that, for a broad range of class-conditional distributions, described by the exponential family, the resulting posterior class probabilities are given by (a logistic (or softmax) transformation acting on a linear function of the feature variables). However, not all choices of class-conditional density give rise to such a simple form for the posterior probabilities (for instance, if the class-conditional densities are modelled using Gaussian mixtures). This suggests that it might be worth exploring other types of discriminative probabilistic model. For the purposes of this chapter, however, we shall return to the two-class case, and again remain within the framework of generalized linear models so that ]

$$p(t=1|a) = f(a) \quad (4.111)$$

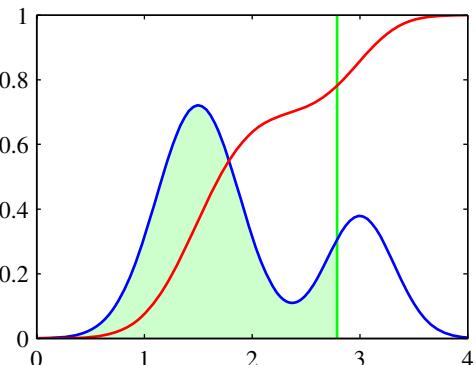
where  $a = w^T \phi$ , and  $f(\cdot)$  is the activation function. ]

[ One way to motivate an alternative choice for the link function is to consider a noisy threshold model, as follows. For each input  $\phi_n$ , we evaluate  $a_n = w^T \phi_n$  and then we set the target value according to ]

$$\begin{cases} t_n = 1 & \text{if } a_n \geq \theta \\ t_n = 0 & \text{otherwise.} \end{cases} \quad (4.112)$$

每次比较，只抽样  
确保并非某一既定值。  
因此,  $t_n$  是实验次  
数越容易时,  $t_n=1$   
出现频率。

**Figure 4.13** Schematic example of a probability density  $p(\theta)$  shown by the blue curve, given in this example by a mixture of two Gaussians, along with its cumulative distribution function  $f(a)$ , shown by the red curve. Note that the value of the blue curve at any point, such as that indicated by the vertical green line, corresponds to the slope of the red curve at the same point. Conversely, the value of the red curve at this point corresponds to the area under the blue curve indicated by the shaded green region. In the stochastic threshold model, the class label takes the value  $t = 1$  if the value of  $a = \mathbf{w}^T \phi$  exceeds a threshold, otherwise it takes the value  $t = 0$ . This is equivalent to an activation function given by the cumulative distribution function  $f(a)$ .



If the value of  $\theta$  is drawn from a probability density  $p(\theta)$ , then the corresponding activation function will be given by the cumulative distribution function

$$f(a) = \int_{-\infty}^a p(\theta) d\theta \quad (4.113)$$

该类判别模型  
的具体实例：  
probit model

as illustrated in Figure 4.13.

As a specific example, suppose that the density  $p(\theta)$  is given by a zero mean, unit variance Gaussian. The corresponding cumulative distribution function is given by

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta|0, 1) d\theta \quad (4.114)$$

which is known as the *probit function*. It has a sigmoidal shape and is compared with the logistic sigmoid function in Figure 4.9. Note that the use of a more general Gaussian distribution does not change the model because this is equivalent to a re-scaling of the linear coefficients  $w$ . Many numerical packages provide for the evaluation of a closely related function defined by

$$\text{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a \exp(-\theta^2) d\theta \quad (4.115)$$

and known as the *erf function* or *error function* (not to be confused with the error function of a machine learning model). It is related to the probit function by

$$\Phi(a) = \frac{1}{2} \left\{ 1 + \frac{1}{\sqrt{2}} \text{erf}(a) \right\}. \quad (4.116)$$

The generalized linear model based on a probit activation function is known as *probit regression*.

参教词学 We can determine the parameters of this model using maximum likelihood, by a straightforward extension of the ideas discussed earlier. In practice, the results found using probit regression tend to be similar to those of logistic regression. We shall,

however, find another use for the probit model when we discuss Bayesian treatments of logistic regression in Section 4.5.

**outliers 的影响** One issue that can occur in practical applications is that of *outliers*, which can arise for instance through errors in measuring the input vector  $\mathbf{x}$  or through mislabelling of the target value  $t$ . Because such points can lie a long way to the wrong side of the ideal decision boundary, they can seriously distort the classifier. Note that the logistic and probit regression models behave differently in this respect because the tails of the logistic sigmoid decay asymptotically like  $\exp(-x)$  for  $x \rightarrow \infty$ , whereas for the probit activation function they decay like  $\exp(-x^2)$ , and so the probit model can be significantly more sensitive to outliers.

**logistic 和 probit 模型均** However, both the logistic and the probit models assume the data is correctly labelled. The effect of mislabelling is easily incorporated into a probabilistic model by introducing a probability  $\epsilon$  that the target value  $t$  has been flipped to the wrong value (Opper and Winther, 2000a), leading to a target value distribution for data point  $\mathbf{x}$  of the form

$$\begin{aligned} p(t|\mathbf{x}) &= (1-\epsilon)\sigma(\mathbf{x}) + \epsilon(1-\sigma(\mathbf{x})) \\ &= \epsilon + (1-2\epsilon)\sigma(\mathbf{x}) \end{aligned} \quad (4.117)$$

where  $\sigma(\mathbf{x})$  is the activation function with input vector  $\mathbf{x}$ . Here  $\epsilon$  may be set in advance, or it may be treated as a hyperparameter whose value is inferred from the data.

### 4.3.6 Canonical link functions

前文介绍的用于回归、分类任务的模型，其损失函数对参数的梯度具有相似的形式

这一现象并不是巧合而是具有一般性的结论

{ For the linear regression model with a Gaussian noise distribution, the error function, corresponding to the negative log likelihood, is given by (3.12). If we take the derivative with respect to the parameter vector  $\mathbf{w}$  of the contribution to the error function from a data point  $n$ , this takes the form of the ‘error’  $y_n - t_n$  times the feature vector  $\phi_n$ , where  $y_n = \mathbf{w}^T \phi_n$ . Similarly, for the combination of the logistic sigmoid activation function and the cross-entropy error function (4.90), and for the softmax activation function with the multiclass cross-entropy error function (4.108), we again obtain this same simple form. } We now show that this is a general result of assuming a conditional distribution for the target variable from the exponential family, along with a corresponding choice for the activation function known as the *canonical link function*.

We again make use of the restricted form (4.84) of exponential family distributions. Note that here we are applying the assumption of exponential family distribution to the target variable  $t$ , in contrast to Section 4.2.4 where we applied it to the input vector  $\mathbf{x}$ . We therefore consider conditional distributions of the target variable of the form

假设 1:  $p(t|\eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\}$ . 是参数,  $s$ ,  $w$  有关,  $\eta(x, w)$  (4.118)

Using the same line of argument as led to the derivation of the result (2.226), we see that the conditional mean of  $t$ , which we denote by  $y$ , is given by

定义:  $y \equiv \mathbb{E}[t|\eta] = -s \frac{d}{d\eta} \ln g(\eta)$ . (4.119)

$\equiv \eta^{-1}(y), \eta(w, x)$

本节书中阐述得有点绕，但总结而言实际上就是：  
假设指数族分布式(4.118)，且其中  $\eta(w, x) = w^T \phi(x)$ ，并记  
 $y \equiv E[t|\eta]$ ，则可推得最终结论式(4.124)。

#### 4.4. The Laplace Approximation

213

Thus  $y$  and  $\eta$  must related, and we denote this relation through  $\eta = \psi(y)$ .

Following Nelder and Wedderburn (1972), we define a *generalized linear model* to be one for which  $y$  is a nonlinear function of a linear combination of the input (or feature) variables so that

**假设2：**  $y = f(w^T \phi)$  即假设  $y = \psi^{-1}(\eta(w, x))$  的函数形式为  $f(w^T \phi(x))$ 。 (4.120)

where  $f(\cdot)$  is known as the *activation function* in the machine learning literature, and  $f^{-1}(\cdot)$  is known as the *link function* in statistics.

Now consider the log likelihood function for this model, which, as a function of  $\eta$ , is given by

$$\ln p(t|\eta, s) = \sum_{n=1}^N \ln p(t_n|\eta, s) = \sum_{n=1}^N \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} + \text{const} \quad (4.121)$$

where we are assuming that all observations share a common scale parameter (which corresponds to the noise variance for a Gaussian distribution for instance) and so  $s$  is independent of  $n$ . The derivative of the log likelihood with respect to the model parameters  $w$  is then given by

**由假设1+假设2 可推得式(4.114)**

$$\begin{aligned} \nabla_w \ln p(t|\eta, s) &= \sum_{n=1}^N \left\{ \frac{d}{d\eta_n} \ln g(\eta_n) + \frac{t_n}{s} \right\} \frac{d\eta_n}{dy_n} \frac{dy_n}{da_n} \nabla a_n \\ &= \sum_{n=1}^N \frac{1}{s} \{t_n - y_n\} \psi'(y_n) f'(a_n) \phi_n \end{aligned} \quad (4.122)$$

where  $a_n = w^T \phi_n$ , and we have used  $y_n = f(a_n)$  together with the result (4.119) for  $E[t|\eta]$ . We now see that there is a considerable simplification if we choose a particular form for the link function  $f^{-1}(y)$  given by 进一步假设函数  $f$ ,  $\psi$  之间的关系

**假设3：**  $f^{-1}(y) = \psi(y)$  又由前文假设  $y = \psi^{-1}(\eta(w, x)) = f(w^T \phi(x))$  (4.123)

which gives  $f(\psi(y)) = y$  and hence  $f'(\psi)\psi'(y) = 1$ . Also, because  $a = f^{-1}(y)$ , we have  $a = \psi$  and hence  $f'(a)\psi'(y) = 1$ . In this case, the gradient of the error function reduces to

**由式(4.122)+假设3 可推得最终结论：**  $\nabla_w E(w) = \frac{1}{s} \sum_{n=1}^N \{y_n - t_n\} \phi_n$  实际上是在假设  $\eta(w, x) = w^T \phi(x)$  (4.124)

For the Gaussian  $s = \beta^{-1}$ , whereas for the logistic model  $s = 1$ .

#### 4.4. The Laplace Approximation

对下文要你做的内容进行说明 [In Section 4.5 we shall discuss the Bayesian treatment of logistic regression. As we shall see, this is more complex than the Bayesian treatment of linear regression models, discussed in Sections 3.3 and 3.5. In particular, we cannot integrate exactly

over the parameter vector  $\mathbf{w}$  since the posterior distribution is no longer Gaussian. It is therefore necessary to introduce some form of approximation. Later in the book we shall consider a range of techniques based on analytical approximations and numerical sampling.]

Here we introduce a simple, but widely used, framework called the Laplace approximation, [that aims to find a Gaussian approximation to a probability density defined over a set of continuous variables.] Consider first the case of a single continuous variable  $z$ , and suppose the distribution  $p(z)$  is defined by

$$p(z) = \frac{1}{Z} f(z) \quad (4.125)$$

where  $Z = \int f(z) dz$  is the normalization coefficient. We shall suppose that the value of  $Z$  is unknown. In the Laplace method the goal is to find a Gaussian approximation  $q(z)$  which is centred on a mode of the distribution  $p(z)$ . [The first step is to find a mode of  $p(z)$ , in other words a point  $z_0$  such that  $p'(z_0) = 0$ , or equivalently

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0. \quad (4.126)$$

A Gaussian distribution has the property that its logarithm is a quadratic function of the variables. We therefore consider a Taylor expansion of  $\ln f(z)$  centred on the mode  $z_0$  so that

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A(z - z_0)^2 \quad (4.127)$$

where

$$A = -\left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}. \quad (4.128)$$

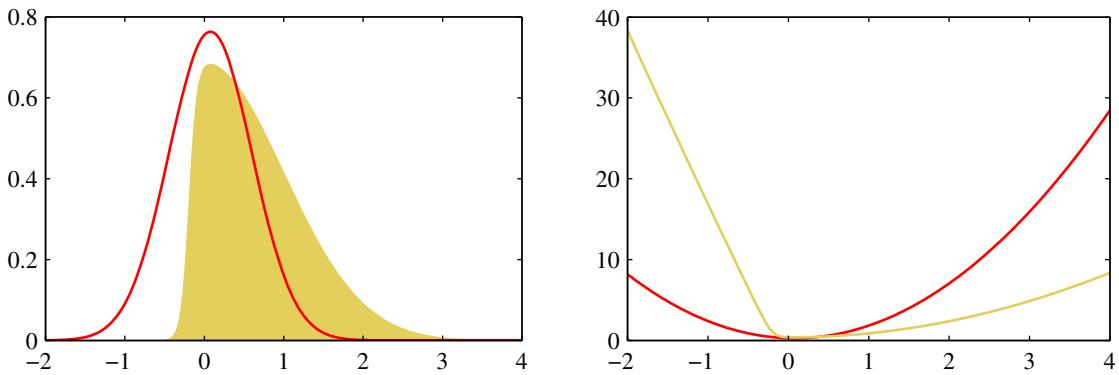
Note that the first-order term in the Taylor expansion does not appear since  $z_0$  is a local maximum of the distribution. Taking the exponential we obtain

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\}. \quad (4.129)$$

We can then obtain a normalized distribution  $q(z)$  by making use of the standard result for the normalization of a Gaussian, so that

$$q(z) = \left( \frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\}. \quad (4.130)$$

The Laplace approximation is illustrated in Figure 4.14. Note that the Gaussian approximation will only be well defined if its precision  $A > 0$ , in other words the stationary point  $z_0$  must be a local maximum, so that the second derivative of  $f(z)$  at the point  $z_0$  is negative.



**Figure 4.14** Illustration of the Laplace approximation applied to the distribution  $p(z) \propto \exp(-z^2/2)\sigma(20z + 4)$  where  $\sigma(z)$  is the logistic sigmoid function defined by  $\sigma(z) = (1 + e^{-z})^{-1}$ . The left plot shows the normalized distribution  $p(z)$  in yellow, together with the Laplace approximation centred on the mode  $z_0$  of  $p(z)$  in red. The right plot shows the negative logarithms of the corresponding curves.

前面介绍的是 Laplace 近似 // We can extend the Laplace method to approximate a distribution  $p(\mathbf{z}) = f(\mathbf{z})/Z$  defined over an  $M$ -dimensional space  $\mathbf{z}$ . At a stationary point  $\mathbf{z}_0$  the gradient  $\nabla f(\mathbf{z})$  will vanish. Expanding around this stationary point we have

$$\ln f(\mathbf{z}) \simeq \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \quad (4.131)$$

where the  $M \times M$  Hessian matrix  $\mathbf{A}$  is defined by

$$\mathbf{A} = -\nabla\nabla \ln f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} \quad (4.132)$$

and  $\nabla$  is the gradient operator. Taking the exponential of both sides we obtain

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\}. \quad (4.133)$$

The distribution  $q(\mathbf{z})$  is proportional to  $f(\mathbf{z})$  and the appropriate normalization coefficient can be found by inspection, using the standard result (2.43) for a normalized multivariate Gaussian, giving

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z} | \mathbf{z}_0, \mathbf{A}^{-1}) \quad (4.134)$$

where  $|\mathbf{A}|$  denotes the determinant of  $\mathbf{A}$ . This Gaussian distribution will be well defined provided its precision matrix, given by  $\mathbf{A}$ , is positive definite, which implies that the stationary point  $\mathbf{z}_0$  must be a local maximum, not a minimum or a saddle point.

In order to apply the Laplace approximation we first need to find the mode  $\mathbf{z}_0$ , and then evaluate the Hessian matrix at that mode. In practice a mode will typically be found by running some form of numerical optimization algorithm (Bishop

and Nabney, 2008). Many of the distributions encountered in practice will be multimodal and so there will be different Laplace approximations according to which mode is being considered. Note that the normalization constant  $Z$  of the true distribution does not need to be known in order to apply the Laplace method. As a result of the central limit theorem, the posterior distribution for a model is expected to become increasingly better approximated by a Gaussian as the number of observed data points is increased, and so we would expect the Laplace approximation to be most useful in situations where the number of data points is relatively large.

*Laplace 近似有2个缺点* One major weakness of the Laplace approximation is that, since it is based on a Gaussian distribution, it is only directly applicable to real variables. In other cases it may be possible to apply the Laplace approximation to a transformation of the variable. For instance if  $0 \leq \tau < \infty$  then we can consider a Laplace approximation of  $\ln \tau$ . The most serious limitation of the Laplace framework, however, is that it is based purely on the aspects of the true distribution at a specific value of the variable, and so can fail to capture important global properties. In Chapter 10 we shall consider alternative approaches which adopt a more global perspective.

#### 4.4.1 Model comparison and BIC

As well as approximating the distribution  $p(\mathbf{z})$  we can also obtain an approximation to the normalization constant  $Z$ . Using the approximation (4.133) we have

$$\begin{aligned} Z &= \int f(\mathbf{z}) d\mathbf{z} \\ &\simeq f(\mathbf{z}_0) \int \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} d\mathbf{z} \\ &= f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \end{aligned} \quad (4.135)$$

where we have noted that the integrand is Gaussian and made use of the standard result (2.43) for a normalized Gaussian distribution. We can use the result (4.135) to obtain an approximation to the model evidence which, as discussed in Section 3.4, plays a central role in Bayesian model comparison.

*model evidence* Consider a data set  $\mathcal{D}$  and a set of models  $\{\mathcal{M}_i\}$  having parameters  $\{\boldsymbol{\theta}_i\}$ . For each model we define a likelihood function  $p(\mathcal{D}|\boldsymbol{\theta}_i, \mathcal{M}_i)$ . If we introduce a prior  $p(\boldsymbol{\theta}_i|\mathcal{M}_i)$  over the parameters, then we are interested in computing the model evidence  $p(\mathcal{D}|\mathcal{M}_i)$  for the various models. From now on we omit the conditioning on  $\mathcal{M}_i$  to keep the notation uncluttered. From Bayes' theorem the model evidence is given by

$$p(\mathcal{D}|\mathcal{D}) = \frac{p(\mathcal{D}, \boldsymbol{\theta})}{p(\mathcal{D})} \underset{f(\boldsymbol{\theta})}{\approx} \frac{1}{2} f(\boldsymbol{\theta}) p(\mathcal{D}) = \int p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (4.136)$$

Identifying  $f(\boldsymbol{\theta}) = p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$  and  $Z = p(\mathcal{D})$ , and applying the result (4.135), we obtain

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \ln p(\boldsymbol{\theta}_{\text{MAP}}) + \underbrace{\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|}_{\text{Occam factor}} \quad (4.137)$$

*Exercise 4.22*

where  $\theta_{\text{MAP}}$  is the value of  $\theta$  at the mode of the posterior distribution, and  $\mathbf{A}$  is the Hessian matrix of second derivatives of the negative log posterior

$$\mathbf{A} = -\nabla\nabla \ln p(\mathcal{D}|\theta_{\text{MAP}})p(\theta_{\text{MAP}}) = -\nabla\nabla \ln p(\theta_{\text{MAP}}|\mathcal{D}). \quad (4.138)$$

The first term on the right hand side of (4.137) represents the log likelihood evaluated using the optimized parameters, while the remaining three terms comprise the ‘Occam factor’ which penalizes model complexity.]

**进一步近似可得 BIC**  
Exercise 4.23

[ If we assume that the Gaussian prior distribution over parameters is broad, and that the Hessian has full rank, then we can approximate (4.137) very roughly using

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\theta_{\text{MAP}}) - \frac{1}{2}M \ln N \quad (4.139)$$

where  $N$  is the number of data points,  $M$  is the number of parameters in  $\theta$  and we have omitted additive constants. This is known as the *Bayesian Information Criterion* (BIC) or the *Schwarz criterion* (Schwarz, 1978). Note that, compared to AIC given by (1.73), this penalizes model complexity more heavily.]

Complexity measures such as AIC and BIC have the virtue of being easy to evaluate, but can also give misleading results. In particular, the assumption that the Hessian matrix has full rank is often not valid since many of the parameters are not ‘well-determined’. We can use the result (4.137) to obtain a more accurate estimate of the model evidence starting from the Laplace approximation, as we illustrate in the context of neural networks in Section 5.7.

### Section 3.5.3

## 4.5. Bayesian Logistic Regression

Bayesian 诊断 (参数  
根据后验分布) 和预  
测无法精确进行  
(Intractable), 因此我们  
借助 Laplace approximation  
进行近似计算。

[ We now turn to a Bayesian treatment of logistic regression. Exact Bayesian inference for logistic regression is intractable. In particular, evaluation of the posterior distribution would require normalization of the product of a prior distribution and a likelihood function that itself comprises a product of logistic sigmoid functions, one for every data point. Evaluation of the predictive distribution is similarly intractable. Here we consider the application of the Laplace approximation to the problem of Bayesian logistic regression (Spiegelhalter and Lauritzen, 1990; MacKay, 1992b). ]

### 4.5.1 Laplace approximation

Recall from Section 4.4 that the Laplace approximation is obtained by finding the mode of the posterior distribution and then fitting a Gaussian centred at that mode. This requires evaluation of the second derivatives of the log posterior, which is equivalent to finding the Hessian matrix.

Because we seek a Gaussian representation for the posterior distribution, it is natural to begin with a Gaussian prior, which we write in the general form

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) \quad (4.140)$$

where  $\mathbf{m}_0$  and  $\mathbf{S}_0$  are fixed hyperparameters. The posterior distribution over  $\mathbf{w}$  is given by

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{w}) \quad (4.141)$$

where  $\mathbf{t} = (t_1, \dots, t_N)^T$ . Taking the log of both sides, and substituting for the prior distribution using (4.140), and for the likelihood function using (4.89), we obtain

$$\begin{aligned} \ln p(\mathbf{w}|\mathbf{t}) &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &\quad + \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + \text{const} \end{aligned} \quad (4.142)$$

对后验分布进行 Laplace 近似

where  $y_n = \sigma(\mathbf{w}^T \phi_n)$ . To obtain a Gaussian approximation to the posterior distribution, we first maximize the posterior distribution to give the MAP (maximum posterior) solution  $\mathbf{w}_{\text{MAP}}$ , which defines the mean of the Gaussian. The covariance is then given by the inverse of the matrix of second derivatives of the negative log likelihood, which takes the form

$$\mathbf{S}_N^{-1} = -\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n) \phi_n \phi_n^T. \quad (4.143)$$

The Gaussian approximation to the posterior distribution therefore takes the form

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{S}_N). \quad (4.144)$$

3. 基于 Laplace approximation 得到  
Having obtained a Gaussian approximation to the posterior distribution, there remains the task of marginalizing with respect to this distribution in order to make predictions.

## 4.5.2 Predictive distribution

[The predictive distribution for class  $\mathcal{C}_1$ , given a new feature vector  $\phi(\mathbf{x})$ , is obtained by marginalizing with respect to the posterior distribution  $p(\mathbf{w}|\mathbf{t})$ , which is itself approximated by a Gaussian distribution  $q(\mathbf{w})$  so that

基于上一节 Laplace approximation 的结果, 进行第一次近似

$$p(\mathcal{C}_1|\phi, \mathbf{t}) = \int p(\mathcal{C}_1|\phi, \mathbf{w}) p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} \quad (4.145)$$

with the corresponding probability for class  $\mathcal{C}_2$  given by  $p(\mathcal{C}_2|\phi, \mathbf{t}) = 1 - p(\mathcal{C}_1|\phi, \mathbf{t})$ . To evaluate the predictive distribution, we first note that the function  $\sigma(\mathbf{w}^T \phi)$  depends on  $\mathbf{w}$  only through its projection onto  $\phi$ . Denoting  $a = \mathbf{w}^T \phi$ , we have

这里很容易让人误认为是  $a = \mathbf{w}^T \phi$ , 但实际上  
我们是积分变量  $a$ , 不等于  $\mathbf{w}^T \phi$   
括号里,  $a$  是积分变量, 而  $\mathbf{w}$  是固定的

$$\sigma(\mathbf{w}^T \phi) = \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da \quad (4.146)$$

两种定义:

where  $\delta(\cdot)$  is the Dirac delta function. From this we obtain

$$\delta(x) = \frac{dH(x)}{dx}, \text{ 其中}$$

$$\delta(x) = \infty, (x \neq 0)$$

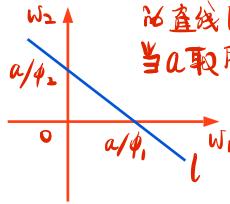
$$\int_{-\infty}^{+\infty} \delta(x) dx = 1$$

$$\delta(\cdot) \text{ 的性质 } \int_{-\infty}^{+\infty} f(x) \delta(x - t_0) dx = f(t_0)$$

$$\int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} = \int \sigma(a) p(a) da \quad (4.147)$$

将式 (4.146) 代入式 (4.145) 后  
变换积分顺序即得式 (4.147)  
和式 (4.148)

对式(4.48)的理解:  
 积分里,  $w$ 是积分变量,  $a/\phi$ 是固定的。当变量  $w$  不满足线性约束  $a - w^T \phi = 0$  时,  $\delta(a - w^T \phi) = 0$ , 则对式(4.48)  
 我们只需计算线性约束  $a - w^T \phi = 0$  成立区域的积分即可。  
 情形为仅: 式(4.48)的积分 <sup>where</sup> 只需在下图中的蓝色直线上进行, 即  $P(a) = \int_L q(w) dw$ 。可以看到, 对不同的  $a$ , 积分  
 直线  $L$  不同, 但它们是平行的, 方向  $w_1$ 。当  $a$  取所有值时,  $\int p(a) da = \int q(w) dw = 1$ 。类似  $P(x_1 | a)$  积分如梯度边缘分布  $P(x_1)$ , 而



$$P(x_1 = a) = \int_{w_1} p(x_1, w_1) dw_1$$

$x_1 = a$   
 积分所在的直线  $L$  的方向向量为  $\{0, 1\}$ ,  $P(a)$  相当于  $q(w)$  的边缘分布。在计

算  $P(a)$  时, 对  $q(w)$  边缘化(即积分)的方向与向量中垂直, 边缘化的位置则由  $a$  给出, 即  $w^T \phi = a$  (类似于  $x_1 = a$ )。又  $q(w)$  为高斯分布, 则其边缘分布  $P(a)$  也是高斯分布。

We can evaluate  $p(a)$  by noting that the delta function imposes a linear constraint on  $w$  and so forms a marginal distribution from the joint distribution  $q(w)$  by integrating out all directions orthogonal to  $\phi$ . Because  $q(w)$  is Gaussian, we know from Section 2.3.2 that the marginal distribution will also be Gaussian. We can evaluate the mean and covariance of this distribution by taking moments, and interchanging the order of integration over  $a$  and  $w$ , so that

$$\mu_a = \mathbb{E}[a] = \int p(a) a da = \int q(w) w^T \phi dw = w_{MAP}^T \phi \quad (4.149)$$

where we have used the result (4.144) for the variational posterior distribution  $q(w)$ . Similarly

$$\sigma_a^2 = \text{var}[a] = \int p(a) \{a^2 - \mathbb{E}[a]^2\} da$$

得式(4.149)  
而此例类似  $\int q(w) \{(w^T \phi)^2 - (\mathbf{m}_N^T \phi)^2\} dw = \phi^T \mathbf{S}_N \phi$ .  
 $\mathbf{m}_N \equiv w_{MAP}$

Note that the distribution of  $a$  takes the same form as the predictive distribution (3.58) for the linear regression model, with the noise variance set to zero. Thus our variational approximation to the predictive distribution becomes

$$p(C_1 | t) = \int \sigma(a) p(a) da = \int \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da. \quad (4.151)$$

这里条件项省略了

This result can also be derived directly by making use of the results for the marginal of a Gaussian distribution given in Section 2.3.2.

The integral over  $a$  represents the convolution of a Gaussian with a logistic sigmoid, and cannot be evaluated analytically. We can, however, obtain a good approximation (Spiegelhalter and Lauritzen, 1990; MacKay, 1992b; Barber and Bishop, 1998a) by making use of the close similarity between the logistic sigmoid function  $\sigma(a)$  defined by (4.59) and the <sup>inverse</sup>probit function  $\Phi(a)$  defined by (4.114). In order to obtain the best approximation to the logistic function we need to re-scale the horizontal axis, so that we approximate  $\sigma(a)$  by  $\Phi(\lambda a)$ . We can find a suitable value of  $\lambda$  by requiring that the two functions have the same slope at the origin, which gives  $\lambda^2 = \pi/8$ . The similarity of the logistic sigmoid and the <sup>inverse</sup>probit function, for this choice of  $\lambda$ , is illustrated in Figure 4.9.

The advantage of using a <sup>inverse</sup>probit function is that its convolution with a Gaussian can be expressed analytically in terms of another <sup>inverse</sup>probit function. Specifically we can show that

$$\int \Phi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right). \quad (4.152)$$

#### Exercise 4.24

对式(4.151)作进一步  
 而近似(用正数近似)  
 积分中的函数  $\sigma(a)$ , 近似  
 后的积分可解析求解。

#### Exercise 4.25

#### Exercise 4.26

$$\int \Phi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right). \quad (4.152)$$

We now apply the approximation  $\sigma(a) \simeq \Phi(\lambda a)$  to the <sup>inverse</sup>probit functions appearing on both sides of this equation, leading to the following approximation for the convolution of a logistic sigmoid with a Gaussian

$$\int \sigma(a)\mathcal{N}(a|\mu, \sigma^2) da \simeq \sigma\left(\kappa(\sigma^2)\mu\right) \quad (4.153)$$

where we have defined

$$\kappa(\sigma^2) = (1 + \pi\sigma^2/8)^{-1/2}. \quad (4.154)$$

Applying this result to (4.151) we obtain the approximate predictive distribution in the form

$$p(C_1|\phi, \mathbf{t}) = \sigma\left(\kappa(\sigma_a^2)\mu_a\right) \quad (4.155)$$

where  $\mu_a$  and  $\sigma_a^2$  are defined by (4.149) and (4.150), respectively, and  $\kappa(\sigma_a^2)$  is defined by (4.154).  $\square$

Note that the decision boundary corresponding to  $p(C_1|\phi, \mathbf{t}) = 0.5$  is given by  $\mu_a = 0$ , which is the same as the decision boundary obtained by using the MAP value for  $\mathbf{w}$ . Thus if the decision criterion is based on minimizing misclassification rate, with equal prior probabilities, then the marginalization over  $\mathbf{w}$  has no effect. However, for more complex decision criteria it will play an important role. Marginalization of the logistic sigmoid model under a Gaussian approximation to the posterior distribution will be illustrated in the context of variational inference in Figure 10.13.

## Exercises

- 4.1** (\*\*) Given a set of data points  $\{\mathbf{x}_n\}$ , we can define the *convex hull* to be the set of all points  $\mathbf{x}$  given by

$$\mathbf{x} = \sum_n \alpha_n \mathbf{x}_n \quad (4.156)$$

where  $\alpha_n \geq 0$  and  $\sum_n \alpha_n = 1$ . Consider a second set of points  $\{\mathbf{y}_n\}$  together with their corresponding convex hull. By definition, the two sets of points will be linearly separable if there exists a vector  $\hat{\mathbf{w}}$  and a scalar  $w_0$  such that  $\hat{\mathbf{w}}^T \mathbf{x}_n + w_0 > 0$  for all  $\mathbf{x}_n$ , and  $\hat{\mathbf{w}}^T \mathbf{y}_n + w_0 < 0$  for all  $\mathbf{y}_n$ . Show that if their convex hulls intersect, the two sets of points cannot be linearly separable, and conversely that if they are linearly separable, their convex hulls do not intersect.

- 4.2** (\*\*) **www** Consider the minimization of a sum-of-squares error function (4.15), and suppose that all of the target vectors in the training set satisfy a linear constraint

$$\mathbf{a}^T \mathbf{t}_n + b = 0 \quad (4.157)$$

where  $\mathbf{t}_n$  corresponds to the  $n^{\text{th}}$  row of the matrix  $\mathbf{T}$  in (4.15). Show that as a consequence of this constraint, the elements of the model prediction  $\mathbf{y}(\mathbf{x})$  given by the least-squares solution (4.17) also satisfy this constraint, so that

$$\mathbf{a}^T \mathbf{y}(\mathbf{x}) + b = 0. \quad (4.158)$$

To do so, assume that one of the basis functions  $\phi_0(\mathbf{x}) = 1$  so that the corresponding parameter  $w_0$  plays the role of a bias.

- 4.3** (\*\*) Extend the result of Exercise 4.2 to show that if multiple linear constraints are satisfied simultaneously by the target vectors, then the same constraints will also be satisfied by the least-squares prediction of a linear model. 4.22
- 4.4** (\*) **www** Show that maximization of the class separation criterion given by (4.23) with respect to  $\mathbf{w}$ , using a Lagrange multiplier to enforce the constraint  $\mathbf{w}^T \mathbf{w} = 1$ , leads to the result that  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ . 4.23
- 4.5** (\*) By making use of (4.20), (4.23), and (4.24), show that the Fisher criterion (4.25) can be written in the form (4.26).
- 4.6** (\*) Using the definitions of the between-class and within-class covariance matrices given by (4.27) and (4.28), respectively, together with (4.34) and (4.36) and the choice of target values described in Section 4.1.5, show that the expression (4.33) that minimizes the sum-of-squares error function can be written in the form (4.37).
- 4.7** (\*) **www** Show that the logistic sigmoid function (4.59) satisfies the property  $\sigma(-a) = 1 - \sigma(a)$  and that its inverse is given by  $\sigma^{-1}(y) = \ln\{y/(1-y)\}$ .
- 4.8** (\*) Using (4.57) and (4.58), derive the result (4.65) for the posterior class probability in the two-class generative model with Gaussian densities, and verify the results (4.66) and (4.67) for the parameters  $\mathbf{w}$  and  $w_0$ .
- 4.9** (\*) **www** Consider a generative classification model for  $K$  classes defined by prior class probabilities  $p(\mathcal{C}_k) = \pi_k$  and general class-conditional densities  $p(\phi|\mathcal{C}_k)$  where  $\phi$  is the input feature vector. Suppose we are given a training data set  $\{\phi_n, \mathbf{t}_n\}$  where  $n = 1, \dots, N$ , and  $\mathbf{t}_n$  is a binary target vector of length  $K$  that uses the 1-of- $K$  coding scheme, so that it has components  $t_{nj} = I_{jk}$  if pattern  $n$  is from class  $\mathcal{C}_k$ . Assuming that the data points are drawn independently from this model, show that the maximum-likelihood solution for the prior probabilities is given by

$$\pi_k = \frac{N_k}{N} \quad (4.159)$$

where  $N_k$  is the number of data points assigned to class  $\mathcal{C}_k$ .

- 4.10** (\*\*) Consider the classification model of Exercise 4.9 and now suppose that the class-conditional densities are given by Gaussian distributions with a shared covariance matrix, so that

$$p(\phi|\mathcal{C}_k) = \mathcal{N}(\phi|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}). \quad (4.160)$$

Show that the maximum likelihood solution for the mean of the Gaussian distribution for class  $\mathcal{C}_k$  is given by

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^{N_k} t_{nk} \phi_n \quad (4.161)$$

which represents the mean of those feature vectors assigned to class  $\mathcal{C}_k$ . Similarly, show that the maximum likelihood solution for the shared covariance matrix is given by

$$\Sigma = \sum_{k=1}^K \frac{N_k}{N} \mathbf{S}_k \quad (4.162)$$

where

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^N t_{nk} (\phi_n - \mu_k)(\phi_n - \mu_k)^T. \quad (4.163)$$

Thus  $\Sigma$  is given by a weighted average of the covariances of the data associated with each class, in which the weighting coefficients are given by the prior probabilities of the classes.

- 4.11** (★) Consider a classification problem with  $K$  classes for which the feature vector  $\phi$  has  $M$  components each of which can take  $L$  discrete states. Let the values of the components be represented by a 1-of- $L$  binary coding scheme. Further suppose that, conditioned on the class  $\mathcal{C}_k$ , the  $M$  components of  $\phi$  are independent, so that the class-conditional density factorizes with respect to the feature vector components. Show that the quantities  $a_k$  given by (4.63), which appear in the argument to the softmax function describing the posterior class probabilities, are linear functions of the components of  $\phi$ . Note that this represents an example of the naive Bayes model which is discussed in Section 8.2.2.
- 4.12** (★) **www** Verify the relation (4.88) for the derivative of the logistic sigmoid function defined by (4.59).
- 4.13** (★) **www** By making use of the result (4.88) for the derivative of the logistic sigmoid, show that the derivative of the error function (4.90) for the logistic regression model is given by (4.91).
- 4.14** (★) Show that for a linearly separable data set, the maximum likelihood solution for the logistic regression model is obtained by finding a vector  $w$  whose decision boundary  $w^T \phi(x) = 0$  separates the classes and then taking the magnitude of  $w$  to infinity.
- 4.15** (★) Show that the Hessian matrix  $H$  for the logistic regression model, given by (4.97), is positive definite. Here  $R$  is a diagonal matrix with elements  $y_n(1 - y_n)$ , and  $y_n$  is the output of the logistic regression model for input vector  $x_n$ . Hence show that the error function is a ~~concave~~ convex function of  $w$  and that it has a unique minimum.
- 4.16** (★) Consider a binary classification problem in which each observation  $x_n$  is known to belong to one of two classes, corresponding to  $t = 0$  and  $t = 1$ , and suppose that the procedure for collecting training data is imperfect, so that training points are sometimes mislabelled. For every data point  $x_n$ , instead of having a value  $t_n$  for the class label, we have instead a value  $\pi_n$  representing the probability that  $t_n = 1$ . Given a probabilistic model  $p(t = 1|\phi)$ , write down the log likelihood function appropriate to such a data set.

- 4.17** (★) **www** Show that the derivatives of the softmax activation function (4.104), where the  $a_k$  are defined by (4.105), are given by (4.106).
- 4.18** (★) Using the result (4.91) for the derivatives of the softmax activation function, show that the gradients of the cross-entropy error (4.108) are given by (4.109). 4.106
- 4.19** (★) **www** Write down expressions for the gradient of the log likelihood, as well as the corresponding Hessian matrix, for the probit regression model defined in Section 4.3.5. These are the quantities that would be required to train such a model using IRLS.
- 4.20** (★★) Show that the Hessian matrix for the multiclass logistic regression problem, defined by (4.110), is positive semidefinite. Note that the full Hessian matrix for this problem is of size  $MK \times MK$ , where  $M$  is the number of parameters and  $K$  is the number of classes. To prove the positive semidefinite property, consider the product  $\mathbf{u}^T \mathbf{H} \mathbf{u}$  where  $\mathbf{u}$  is an arbitrary vector of length  $MK$ , and then apply Jensen's inequality. inverse
- 4.21** (★) Show that the probit function (4.114) and the erf function (4.115) are related by (4.116).
- 4.22** (★) Using the result (4.135), derive the expression (4.137) for the log model evidence under the Laplace approximation.
- 4.23** (★★) **www** In this exercise, we derive the BIC result (4.139) starting from the Laplace approximation to the model evidence given by (4.137). Show that if the prior over parameters is Gaussian of the form  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}, \mathbf{V}_0)$ , the log model evidence under the Laplace approximation takes the form

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D} | \boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2} (\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}) - \frac{1}{2} \ln |\mathbf{H}| + \text{const}$$
negative

where  $\mathbf{H}$  is the matrix of second derivatives of the log likelihood  $\ln p(\mathcal{D} | \boldsymbol{\theta})$  evaluated at  $\boldsymbol{\theta}_{\text{MAP}}$ . Now assume that the prior is broad so that  $\mathbf{V}_0^{-1}$  is small and the second term on the right-hand side above can be neglected. Furthermore, consider the case of independent, identically distributed data so that  $\mathbf{H}$  is the sum of terms one for each data point. Show that the log model evidence can then be written approximately in the form of the BIC expression (4.139).

- 4.24** (★★) Use the results from Section 2.3.2 to derive the result (4.151) for the marginalization of the logistic regression model with respect to a Gaussian posterior distribution over the parameters  $\mathbf{w}$ .
- 4.25** (★★) Suppose we wish to approximate the logistic sigmoid  $\sigma(a)$  defined by (4.59) by a scaled probit function  $\Phi(\lambda a)$ , where  $\Phi(a)$  is defined by (4.114). Show that if  $\lambda$  is chosen so that the derivatives of the two functions are equal at  $a = 0$ , then  $\lambda^2 = \pi/8$ . inverse

inverse  
F

- 4.26** (\*\*) In this exercise, we prove the relation (4.152) for the convolution of a probit function with a Gaussian distribution. To do this, show that the derivative of the left-hand side with respect to  $\mu$  is equal to the derivative of the right-hand side, and then integrate both sides with respect to  $\mu$  and then show that the constant of integration vanishes. Note that before differentiating the left-hand side, it is convenient first to introduce a change of variable given by  $a = \mu + \sigma z$  so that the integral over  $a$  is replaced by an integral over  $z$ . When we differentiate the left-hand side of the relation (4.152), we will then obtain a Gaussian integral over  $z$  that can be evaluated analytically.

---

# 5

# Neural Networks

In Chapters 3 and 4 we considered models for regression and classification that comprised linear combinations of fixed basis functions. We saw that such models have useful analytical and computational properties but that their practical applicability was limited by the curse of dimensionality. In order to apply such models to large-scale problems, it is necessary to adapt the basis functions to the data.

Support vector machines (SVMs), discussed in Chapter 7, address this by first defining basis functions that are centred on the training data points and then selecting a subset of these during training. One advantage of SVMs is that, although the training involves nonlinear optimization, the objective function is convex, and so the solution of the optimization problem is relatively straightforward. The number of basis functions in the resulting models is generally much smaller than the number of training points, although it is often still relatively large and typically increases with the size of the training set. The relevance vector machine, discussed in Section 7.2, also chooses a subset from a fixed set of basis functions and typically results in much

sparser models. Unlike the SVM it also produces probabilistic outputs, although this is at the expense of a nonconvex optimization during training.

An alternative approach is to fix the number of basis functions in advance but allow them to be adaptive, in other words to use parametric forms for the basis functions in which the parameter values are adapted during training. The most successful model of this type in the context of pattern recognition is the feed-forward neural network, also known as the *multilayer perceptron*, discussed in this chapter. In fact, ‘multilayer perceptron’ is really a misnomer, because the model comprises multiple layers of logistic regression models (with continuous nonlinearities) rather than multiple perceptrons (with discontinuous nonlinearities). For many applications, the resulting model can be significantly more compact, and hence faster to evaluate, than a support vector machine having the same generalization performance. The price to be paid for this compactness, as with the relevance vector machine, is that the likelihood function, which forms the basis for network training, is no longer a convex function of the model parameters. In practice, however, it is often worth investing substantial computational resources during the training phase in order to obtain a compact model that is fast at processing new data.

The term ‘neural network’ has its origins in attempts to find mathematical representations of information processing in biological systems (McCulloch and Pitts, 1943; Widrow and Hoff, 1960; Rosenblatt, 1962; Rumelhart *et al.*, 1986). Indeed, it has been used very broadly to cover a wide range of different models, many of which have been the subject of exaggerated claims regarding their biological plausibility. From the perspective of practical applications of pattern recognition, however, biological realism would impose entirely unnecessary constraints. Our focus in this chapter is therefore on neural networks as efficient models for statistical pattern recognition. In particular, we shall restrict our attention to the specific class of neural networks that have proven to be of greatest practical value, namely the multilayer perceptron.

**本章要讨论的内容** We begin by considering the functional form of the network model, including the specific parameterization of the basis functions, and we then discuss the problem of determining the network parameters within a maximum likelihood framework, which involves the solution of a nonlinear optimization problem. This requires the evaluation of derivatives of the log likelihood function with respect to the network parameters, and we shall see how these can be obtained efficiently using the technique of *error backpropagation*. We shall also show how the backpropagation framework can be extended to allow other derivatives to be evaluated, such as the Jacobian and Hessian matrices. Next we discuss various approaches to regularization of neural network training and the relationships between them. We also consider some extensions to the neural network model, and in particular we describe a general framework for modelling conditional probability distributions known as *mixture density networks*. Finally, we discuss the use of Bayesian treatments of neural networks. Additional background on neural network models can be found in Bishop (1995a).

## 5.1. Feed-forward Network Functions

The linear models for regression and classification discussed in Chapters 3 and 4, respectively, are based on linear combinations of fixed nonlinear basis functions  $\phi_j(\mathbf{x})$  and take the form

$$y(\mathbf{x}, \mathbf{w}) = f \left( \sum_{j=1}^M w_j \phi_j(\mathbf{x}) \right) \quad (5.1)$$

where  $f(\cdot)$  is a nonlinear activation function in the case of classification and is the identity in the case of regression. Our goal is to extend this model by making the basis functions  $\phi_j(\mathbf{x})$  depend on parameters and then to allow these parameters to be adjusted, along with the coefficients  $\{w_j\}$ , during training. There are, of course, many ways to construct parametric nonlinear basis functions. Neural networks use basis functions that follow the same form as (5.1), so that each basis function is itself a nonlinear function of a linear combination of the inputs, where the coefficients in the linear combination are adaptive parameters.

This leads to the basic neural network model, which can be described as a series of functional transformations. First we construct  $M$  linear combinations of the input variables  $x_1, \dots, x_D$  in the form

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \quad (5.2)$$

where  $j = 1, \dots, M$ , and the superscript (1) indicates that the corresponding parameters are in the first ‘layer’ of the network. We shall refer to the parameters  $w_{ji}^{(1)}$  as **weights** and the parameters  $w_{j0}^{(1)}$  as **biases**, following the nomenclature of Chapter 3. The quantities  $a_j$  are known as **activations**. Each of them is then transformed using a differentiable, nonlinear *activation function*  $h(\cdot)$  to give

$$z_j = h(a_j). \quad (5.3)$$

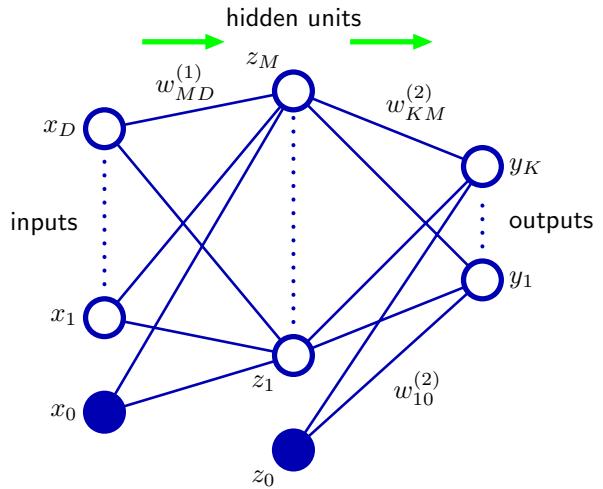
These quantities correspond to the outputs of the basis functions in (5.1) that, in the context of neural networks, are called *hidden units*. The nonlinear functions  $h(\cdot)$  are generally chosen to be sigmoidal functions such as the logistic sigmoid or the ‘tanh’ function. Following (5.1), these values are again linearly combined to give *output unit activations*

$$a_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + w_{k0}^{(2)} \quad (5.4)$$

where  $k = 1, \dots, K$ , and  $K$  is the total number of outputs. This transformation corresponds to the second layer of the network, and again the  $w_{k0}^{(2)}$  are bias parameters. Finally, the output unit activations are transformed using an appropriate activation function to give a set of network outputs  $y_k$ . The choice of activation function is determined by the nature of the data and the assumed distribution of target variables

### Exercise 5.1

**Figure 5.1** Network diagram for the two-layer neural network corresponding to (5.7). The input, hidden, and output variables are represented by nodes, and the weight parameters are represented by links between the nodes, in which the bias parameters are denoted by links coming from additional input and hidden variables  $x_0$  and  $z_0$ . Arrows denote the direction of information flow through the network during forward propagation.



and follows the same considerations as for linear models discussed in Chapters 3 and 4. Thus for standard regression problems, the activation function is the identity so that  $y_k = a_k$ . Similarly, for multiple binary classification problems, each output unit activation is transformed using a logistic sigmoid function so that

$$y_k = \sigma(a_k) \quad (5.5)$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}. \quad (5.6)$$

Finally, for multiclass problems, a softmax activation function of the form (4.62) is used. The choice of output unit activation function is discussed in detail in Section 5.2.

We can combine these various stages to give the overall network function that, for sigmoidal output unit activation functions, takes the form

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \sum_{j=1}^M w_{kj}^{(2)} h \left( \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right) \quad (5.7)$$

where the set of all weight and bias parameters have been grouped together into a vector  $\mathbf{w}$ . Thus the neural network model is simply a nonlinear function from a set of input variables  $\{x_i\}$  to a set of output variables  $\{y_k\}$  controlled by a vector  $\mathbf{w}$  of adjustable parameters.

This function can be represented in the form of a network diagram as shown in Figure 5.1. The process of evaluating (5.7) can then be interpreted as a *forward propagation* of information through the network. It should be emphasized that these diagrams do not represent probabilistic graphical models of the kind to be considered in Chapter 8 because the internal nodes represent deterministic variables rather than stochastic ones. For this reason, we have adopted a slightly different graphical

notation for the two kinds of model. We shall see later how to give a probabilistic interpretation to a neural network.

As discussed in Section 3.1, the bias parameters in (5.2) can be absorbed into the set of weight parameters by defining an additional input variable  $x_0$  whose value is clamped at  $x_0 = 1$ , so that (5.2) takes the form

$$a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i. \quad (5.8)$$

We can similarly absorb the second-layer biases into the second-layer weights, so that the overall network function becomes

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \sum_{j=0}^M w_{kj}^{(2)} h \left( \sum_{i=0}^D w_{ji}^{(1)} x_i \right) \right). \quad (5.9)$$

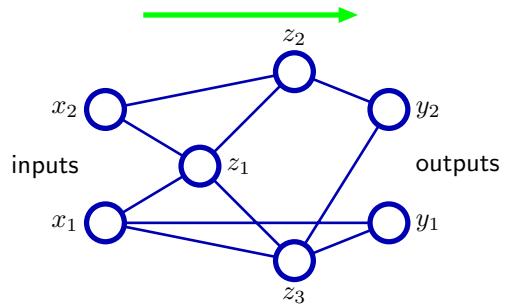
As can be seen from Figure 5.1, the neural network model comprises two stages of processing, each of which resembles the perceptron model of Section 4.1.7, and for this reason the neural network is also known as the *multilayer perceptron*, or MLP. A key difference compared to the perceptron, however, is that the neural network uses continuous sigmoidal nonlinearities in the hidden units, whereas the perceptron uses step-function nonlinearities. This means that the neural network function is differentiable with respect to the network parameters, and this property will play a central role in network training.

If the activation functions of all the hidden units in a network are taken to be linear, then for any such network we can always find an equivalent network without hidden units. This follows from the fact that the composition of successive linear transformations is itself a linear transformation. However, if the number of hidden units is smaller than either the number of input or output units, then the transformations that the network can generate are not the most general possible linear transformations from inputs to outputs because information is lost in the dimensionality reduction at the hidden units. In Section 12.4.2, we show that networks of linear units give rise to principal component analysis. In general, however, there is little interest in multilayer networks of linear units.

The network architecture shown in Figure 5.1 is the most commonly used one in practice. However, it is easily generalized, for instance by considering additional layers of processing each consisting of a weighted linear combination of the form (5.4) followed by an element-wise transformation using a nonlinear activation function. Note that there is some confusion in the literature regarding the terminology for counting the number of layers in such networks. Thus the network in Figure 5.1 may be described as a 3-layer network (which counts the number of layers of units, and treats the inputs as units) or sometimes as a single-hidden-layer network (which counts the number of layers of hidden units). We recommend a terminology in which Figure 5.1 is called a two-layer network, because it is the number of layers of adaptive weights that is important for determining the network properties.

Another generalization of the network architecture is to include *skip-layer* connections, each of which is associated with a corresponding adaptive parameter. For

**Figure 5.2** Example of a neural network having a general feed-forward topology. Note that each hidden and output unit has an associated bias parameter (omitted for clarity).



instance, in a two-layer network these would go directly from inputs to outputs. In principle, a network with sigmoidal hidden units can always mimic skip layer connections (for bounded input values) by using a sufficiently small first-layer weight that, over its operating range, the hidden unit is effectively linear, and then compensating with a large weight value from the hidden unit to the output. In practice, however, it may be advantageous to include skip-layer connections explicitly.

Furthermore, the network can be sparse, with not all possible connections within a layer being present. We shall see an example of a sparse network architecture when we consider convolutional neural networks in Section 5.5.6.

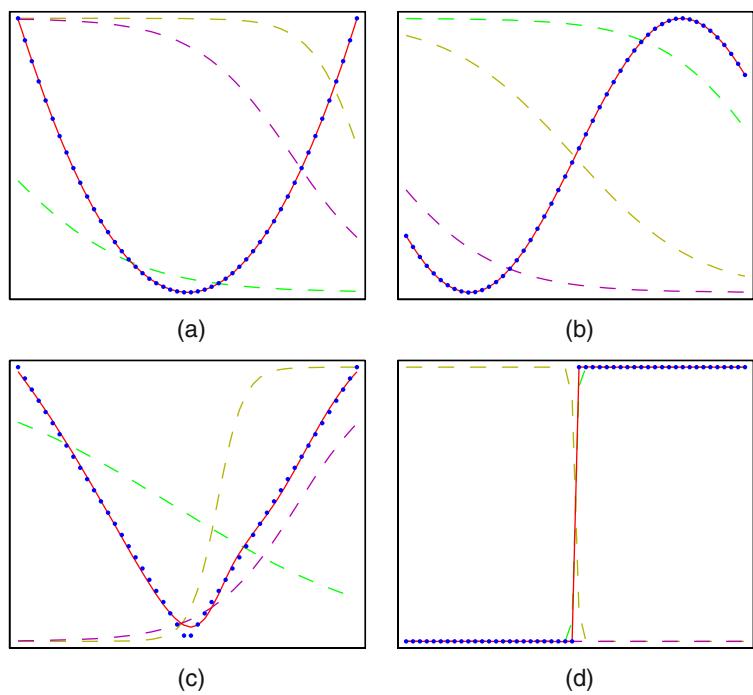
Because there is a direct correspondence between a network diagram and its mathematical function, we can develop more general network mappings by considering more complex network diagrams. However, these must be restricted to a *feed-forward* architecture, in other words to one having no closed directed cycles, to ensure that the outputs are deterministic functions of the inputs. This is illustrated with a simple example in Figure 5.2. Each (hidden or output) unit in such a network computes a function given by

$$z_k = h \left( \sum_j w_{kj} z_j \right) \quad (5.10)$$

where the sum runs over all units that send connections to unit  $k$  (and a bias parameter is included in the summation). For a given set of values applied to the inputs of the network, successive application of (5.10) allows the activations of all units in the network to be evaluated including those of the output units.

The approximation properties of feed-forward networks have been widely studied (Funahashi, 1989; Cybenko, 1989; Hornik *et al.*, 1989; Stinchcombe and White, 1989; Cotter, 1990; Ito, 1991; Hornik, 1991; Kreinovich, 1991; Ripley, 1996) and found to be very general. Neural networks are therefore said to be *universal approximators*. For example, a two-layer network with linear outputs can uniformly approximate any continuous function on a compact input domain to arbitrary accuracy provided the network has a sufficiently large number of hidden units. This result holds for a wide range of hidden unit activation functions, but excluding polynomials. Although such theorems are reassuring, the key problem is how to find suitable parameter values given a set of training data, and in later sections of this chapter we

**Figure 5.3** Illustration of the capability of a multilayer perceptron to approximate four different functions comprising (a)  $f(x) = x^2$ , (b)  $f(x) = \sin(x)$ , (c),  $f(x) = |x|$ , and (d)  $f(x) = H(x)$  where  $H(x)$  is the Heaviside step function. In each case,  $N = 50$  data points, shown as blue dots, have been sampled uniformly in  $x$  over the interval  $(-1, 1)$  and the corresponding values of  $f(x)$  evaluated. These data points are then used to train a two-layer network having 3 hidden units with ‘tanh’ activation functions and linear output units. The resulting network functions are shown by the red curves, and the outputs of the three hidden units are shown by the three dashed curves.



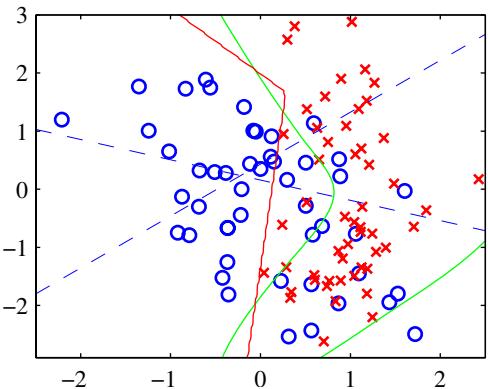
will show that there exist effective solutions to this problem based on both maximum likelihood and Bayesian approaches.

The capability of a two-layer network to model a broad range of functions is illustrated in Figure 5.3. This figure also shows how individual hidden units work collaboratively to approximate the final function. The role of hidden units in a simple classification problem is illustrated in Figure 5.4 using the synthetic classification data set described in Appendix A.

### 5.1.1 Weight-space symmetries

One property of feed-forward networks, which will play a role when we consider Bayesian model comparison, is that multiple distinct choices for the weight vector  $w$  can all give rise to the same mapping function from inputs to outputs (Chen *et al.*, 1993). Consider a two-layer network of the form shown in Figure 5.1 with  $M$  hidden units having ‘tanh’ activation functions and full connectivity in both layers. If we change the sign of all of the weights and the bias feeding into a particular hidden unit, then, for a given input pattern, the sign of the activation of the hidden unit will be reversed, because ‘tanh’ is an odd function, so that  $\tanh(-a) = -\tanh(a)$ . This transformation can be exactly compensated by changing the sign of all of the weights leading out of that hidden unit. Thus, by changing the signs of a particular group of weights (and a bias), the input–output mapping function represented by the network is unchanged, and so we have found two different weight vectors that give rise to the same mapping function. For  $M$  hidden units, there will be  $M$  such ‘sign-flip’

**Figure 5.4** Example of the solution of a simple two-class classification problem involving synthetic data using a neural network having two inputs, two hidden units with ‘tanh’ activation functions, and a single output having a logistic sigmoid activation function. The dashed blue lines show the  $z = 0.5$  contours for each of the hidden units, and the red line shows the  $y = 0.5$  decision surface for the network. For comparison, the green line denotes the optimal decision boundary computed from the distributions used to generate the data.



symmetries, and thus any given weight vector will be one of a set  $2^M$  equivalent weight vectors .

Similarly, imagine that we interchange the values of all of the weights (and the bias) leading both into and out of a particular hidden unit with the corresponding values of the weights (and bias) associated with a different hidden unit. Again, this clearly leaves the network input–output mapping function unchanged, but it corresponds to a different choice of weight vector. For  $M$  hidden units, any given weight vector will belong to a set of  $M!$  equivalent weight vectors associated with this interchange symmetry, corresponding to the  $M!$  different orderings of the hidden units. The network will therefore have an overall weight-space symmetry factor of  $M!2^M$ . For networks with more than two layers of weights, the total level of symmetry will be given by the product of such factors, one for each layer of hidden units.

It turns out that these factors account for all of the symmetries in weight space (except for possible accidental symmetries due to specific choices for the weight values). Furthermore, the existence of these symmetries is not a particular property of the ‘tanh’ function but applies to a wide range of activation functions (Kúrková and Kainen, 1994). In many cases, these symmetries in weight space are of little practical consequence, although in Section 5.7 we shall encounter a situation in which we need to take them into account.

## 5.2. Network Training

So far, we have viewed neural networks as a general class of parametric nonlinear functions from a vector  $\mathbf{x}$  of input variables to a vector  $\mathbf{y}$  of output variables. A simple approach to the problem of determining the network parameters is to make an analogy with the discussion of polynomial curve fitting in Section 1.1, and therefore to minimize a sum-of-squares error function. Given a training set comprising a set of input vectors  $\{\mathbf{x}_n\}$ , where  $n = 1, \dots, N$ , together with a corresponding set of

target vectors  $\{\mathbf{t}_n\}$ , we minimize the error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|y(\mathbf{x}_n, \mathbf{w}) - t_n\|^2. \quad (5.11)$$

However, we can provide a much more general view of network training by first giving a probabilistic interpretation to the network outputs. We have already seen many advantages of using probabilistic predictions in Section 1.5.4. Here it will also provide us with a clearer motivation both for the choice of output unit nonlinearity and the choice of error function.

回归问题的  
损失函数

We start by discussing regression problems, and for the moment we consider a single target variable  $t$  that can take any real value. Following the discussions in Section 1.2.5 and 3.1, we assume that  $t$  has a Gaussian distribution with an  $\mathbf{x}$ -dependent mean, which is given by the output of the neural network, so that

$$p(t|\mathbf{x}, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \quad (5.12)$$

where  $\beta$  is the precision (inverse variance) of the Gaussian noise. Of course this is a somewhat restrictive assumption, and in Section 5.6 we shall see how to extend this approach to allow for more general conditional distributions. For the conditional distribution given by (5.12), it is sufficient to take the output unit activation function to be the identity, because such a network can approximate any continuous function from  $\mathbf{x}$  to  $y$ . Given a data set of  $N$  independent, identically distributed observations  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , along with corresponding target values  $\mathbf{t} = \{t_1, \dots, t_N\}$ , we can construct the corresponding likelihood function

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \beta).$$

Taking the negative logarithm, we obtain the error function

$$\frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln(2\pi) \quad (5.13)$$

which can be used to learn the parameters  $\mathbf{w}$  and  $\beta$ . In Section 5.7, we shall discuss the Bayesian treatment of neural networks, while here we consider a maximum likelihood approach. Note that in the neural networks literature, it is usual to consider the minimization of an error function rather than the maximization of the (log) likelihood, and so here we shall follow this convention. Consider first the determination of  $\mathbf{w}$ . Maximizing the likelihood function is equivalent to minimizing the sum-of-squares error function given by

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 \quad (5.14)$$

where we have discarded additive and multiplicative constants. The value of  $\mathbf{w}$  found by minimizing  $E(\mathbf{w})$  will be denoted  $\mathbf{w}_{\text{ML}}$  because it corresponds to the maximum likelihood solution. In practice, the nonlinearity of the network function  $y(\mathbf{x}_n, \mathbf{w})$  causes the error  $E(\mathbf{w})$  to be nonconvex, and so in practice local maxima of the likelihood may be found, corresponding to local minima of the error function, as discussed in Section 5.2.1.

Having found  $\mathbf{w}_{\text{ML}}$ , the value of  $\beta$  can be found by minimizing the negative log likelihood to give

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}_{\text{ML}}) - t_n\}^2. \quad (5.15)$$

Note that this can be evaluated once the iterative optimization required to find  $\mathbf{w}_{\text{ML}}$  is completed. If we have multiple target variables, and we assume that they are independent conditional on  $\mathbf{x}$  and  $\mathbf{w}$  with shared noise precision  $\beta$ , then the conditional distribution of the target values is given by

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{t}|y(\mathbf{x}, \mathbf{w}), \beta^{-1}\mathbf{I}). \quad (5.16)$$

Following the same argument as for a single target variable, we see that the maximum likelihood weights are determined by minimizing the sum-of-squares error function (5.11). The noise precision is then given by

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{NK} \sum_{n=1}^N \|y(\mathbf{x}_n, \mathbf{w}_{\text{ML}}) - \mathbf{t}_n\|^2 \quad (5.17)$$

where  $K$  is the number of target variables. The assumption of independence can be dropped at the expense of a slightly more complex optimization problem.

### Exercise 5.2

### Exercise 5.3

Recall from Section 4.3.6 that there is a natural pairing of the error function (given by the negative log likelihood) and the output unit activation function. In the regression case, we can view the network as having an output activation function that is the identity, so that  $y_k = a_k$ . The corresponding sum-of-squares error function has the property

$$\frac{\partial E}{\partial a_k} = y_k - t_k \quad (5.18)$$

分类问题  
损失函数 which we shall make use of when discussing error backpropagation in Section 5.3.] Now consider the case of binary classification in which we have a single target variable  $t$  such that  $t = 1$  denotes class  $C_1$  and  $t = 0$  denotes class  $C_2$ . Following the discussion of canonical link functions in Section 4.3.6, we consider a network having a single output whose activation function is a logistic sigmoid

$$y = \sigma(a) \equiv \frac{1}{1 + \exp(-a)} \quad (5.19)$$

so that  $0 \leq y(\mathbf{x}, \mathbf{w}) \leq 1$ . We can interpret  $y(\mathbf{x}, \mathbf{w})$  as the conditional probability  $p(C_1|\mathbf{x})$ , with  $p(C_2|\mathbf{x})$  given by  $1 - y(\mathbf{x}, \mathbf{w})$ . The conditional distribution of targets given inputs is then a Bernoulli distribution of the form

$$p(t|\mathbf{x}, \mathbf{w}) = y(\mathbf{x}, \mathbf{w})^t \{1 - y(\mathbf{x}, \mathbf{w})\}^{1-t}. \quad (5.20)$$

If we consider a training set of independent observations, then the error function, which is given by the negative log likelihood, is then a *cross-entropy* error function of the form

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (5.21)$$

where  $y_n$  denotes  $y(\mathbf{x}_n, \mathbf{w})$ . Note that there is no analogue of the noise precision  $\beta$  because the target values are assumed to be correctly labelled. However, the model is easily extended to allow for labelling errors. Simard *et al.* (2003) found that using the cross-entropy error function instead of the sum-of-squares for a classification problem leads to faster training as well as improved generalization.

If we have  $K$  separate binary classifications to perform, then we can use a network having  $K$  outputs each of which has a logistic sigmoid activation function. Associated with each output is a binary class label  $t_k \in \{0, 1\}$ , where  $k = 1, \dots, K$ . If we assume that the class labels are independent, given the input vector, then the conditional distribution of the targets is

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \prod_{k=1}^K y_k(\mathbf{x}, \mathbf{w})^{t_k} [1 - y_k(\mathbf{x}, \mathbf{w})]^{1-t_k}. \quad (5.22)$$

Taking the negative logarithm of the corresponding likelihood function then gives the following error function

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K \{t_{nk} \ln y_{nk} + (1 - t_{nk}) \ln(1 - y_{nk})\} \quad (5.23)$$

where  $y_{nk}$  denotes  $y_k(\mathbf{x}_n, \mathbf{w})$ . Again, the derivative of the error function with respect to the activation for a particular output unit takes the form (5.18) just as in the regression case.

It is interesting to contrast the neural network solution to this problem with the corresponding approach based on a linear classification model of the kind discussed in Chapter 4. Suppose that we are using a standard two-layer network of the kind shown in Figure 5.1. We see that the weight parameters in the first layer of the network are shared between the various outputs, whereas in the linear model each classification problem is solved independently. The first layer of the network can be viewed as performing a nonlinear feature extraction, and the sharing of features between the different outputs can save on computation and can also lead to improved generalization.

Finally, we consider the standard multiclass classification problem in which each input is assigned to one of  $K$  mutually exclusive classes. The binary target variables  $t_k \in \{0, 1\}$  have a 1-of- $K$  coding scheme indicating the class, and the network outputs are interpreted as  $y_k(\mathbf{x}, \mathbf{w}) = p(t_k = 1|\mathbf{x})$ , leading to the following error function

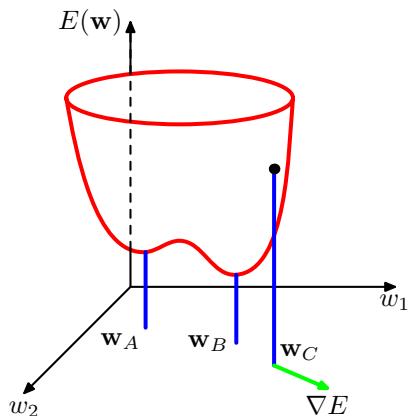
$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K \frac{t_{nk}}{1 - t_{nk}} \ln y_k(\mathbf{x}_n, \mathbf{w}). \quad (5.24)$$

### Exercise 5.4

### Exercise 5.5

### Exercise 5.6

**Figure 5.5** Geometrical view of the error function  $E(\mathbf{w})$  as a surface sitting over weight space. Point  $\mathbf{w}_A$  is a local minimum and  $\mathbf{w}_B$  is the global minimum. At any point  $\mathbf{w}_C$ , the local gradient of the error surface is given by the vector  $\nabla E$ .



Following the discussion of Section 4.3.4, we see that the output unit activation function, which corresponds to the canonical link, is given by the softmax function

$$y_k(\mathbf{x}, \mathbf{w}) = \frac{\exp(a_k(\mathbf{x}, \mathbf{w}))}{\sum_j \exp(a_j(\mathbf{x}, \mathbf{w}))} \quad (5.25)$$

which satisfies  $0 \leq y_k \leq 1$  and  $\sum_k y_k = 1$ . Note that the  $y_k(\mathbf{x}, \mathbf{w})$  are unchanged if a constant is added to all of the  $a_k(\mathbf{x}, \mathbf{w})$ , causing the error function to be constant for some directions in weight space. This degeneracy is removed if an appropriate regularization term (Section 5.5) is added to the error function.

*Exercise 5.7*

Once again, the derivative of the error function with respect to the activation for a particular output unit takes the familiar form (5.18).

In summary, there is a natural choice of both output unit activation function and matching error function, according to the type of problem being solved. For regression we use linear outputs and a sum-of-squares error, for (multiple independent) binary classifications we use logistic sigmoid outputs and a cross-entropy error function, and for multiclass classification we use softmax outputs with the corresponding multiclass cross-entropy error function. For classification problems involving two classes, we can use a single logistic sigmoid output, or alternatively we can use a network with two outputs having a softmax output activation function.

### 5.2.1 Parameter optimization

We turn next to the task of finding a weight vector  $\mathbf{w}$  which minimizes the chosen function  $E(\mathbf{w})$ . At this point, it is useful to have a geometrical picture of the error function, which we can view as a surface sitting over weight space as shown in Figure 5.5. First note that if we make a small step in weight space from  $\mathbf{w}$  to  $\mathbf{w} + \delta\mathbf{w}$  then the change in the error function is  $\delta E \simeq \delta\mathbf{w}^T \nabla E(\mathbf{w})$ , where the vector  $\nabla E(\mathbf{w})$  points in the direction of greatest rate of increase of the error function. Because the error  $E(\mathbf{w})$  is a smooth continuous function of  $\mathbf{w}$ , its smallest value will occur at a

point in weight space such that the gradient of the error function vanishes, so that

$$\nabla E(\mathbf{w}) = 0 \quad (5.26)$$

as otherwise we could make a small step in the direction of  $-\nabla E(\mathbf{w})$  and thereby further reduce the error. Points at which the gradient vanishes are called stationary points, and may be further classified into minima, maxima, and saddle points.

Our goal is to find a vector  $\mathbf{w}$  such that  $E(\mathbf{w})$  takes its smallest value. However, the error function typically has a highly nonlinear dependence on the weights and bias parameters, and so there will be many points in weight space at which the gradient vanishes (or is numerically very small). Indeed, from the discussion in Section 5.1.1 we see that for any point  $\mathbf{w}$  that is a local minimum, there will be other points in weight space that are equivalent minima. For instance, in a two-layer network of the kind shown in Figure 5.1, with  $M$  hidden units, each point in weight space is a member of a family of  $M!2^M$  equivalent points.

### Section 5.1.1

Furthermore, there will typically be multiple inequivalent stationary points and in particular multiple inequivalent minima. A minimum that corresponds to the smallest value of the error function for any weight vector is said to be a *global minimum*. Any other minima corresponding to higher values of the error function are said to be *local minima*. For a successful application of neural networks, it may not be necessary to find the global minimum (and in general it will not be known whether the global minimum has been found) but it may be necessary to compare several local minima in order to find a sufficiently good solution.

Because there is clearly no hope of finding an analytical solution to the equation  $\nabla E(\mathbf{w}) = 0$  we resort to iterative numerical procedures. The optimization of continuous nonlinear functions is a widely studied problem and there exists an extensive literature on how to solve it efficiently. Most techniques involve choosing some initial value  $\mathbf{w}^{(0)}$  for the weight vector and then moving through weight space in a succession of steps of the form

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \Delta\mathbf{w}^{(\tau)} \quad (5.27)$$

where  $\tau$  labels the iteration step. Different algorithms involve different choices for the weight vector update  $\Delta\mathbf{w}^{(\tau)}$ . Many algorithms make use of gradient information and therefore require that, after each update, the value of  $\nabla E(\mathbf{w})$  is evaluated at the new weight vector  $\mathbf{w}^{(\tau+1)}$ . In order to understand the importance of gradient information, it is useful to consider a local approximation to the error function based on a Taylor expansion.

## 5.2.2 Local quadratic approximation

Insight into the optimization problem, and into the various techniques for solving it, can be obtained by considering a local quadratic approximation to the error function.

Consider the Taylor expansion of  $E(\mathbf{w})$  around some point  $\hat{\mathbf{w}}$  in weight space

$$E(\mathbf{w}) \simeq E(\hat{\mathbf{w}}) + (\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{b} + \frac{1}{2} (\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{H} (\mathbf{w} - \hat{\mathbf{w}}) \quad (5.28)$$

where cubic and higher terms have been omitted. Here  $\mathbf{b}$  is defined to be the gradient of  $E$  evaluated at  $\hat{\mathbf{w}}$

$$\mathbf{b} \equiv \nabla E|_{\mathbf{w}=\hat{\mathbf{w}}} \quad (5.29)$$

and the Hessian matrix  $\mathbf{H} = \nabla \nabla E$  has elements

$$(\mathbf{H})_{ij} \equiv \left. \frac{\partial E}{\partial w_i \partial w_j} \right|_{\mathbf{w}=\hat{\mathbf{w}}}. \quad (5.30)$$

From (5.28), the corresponding local approximation to the gradient is given by

$$\nabla E \simeq \mathbf{b} + \mathbf{H}(\mathbf{w} - \hat{\mathbf{w}}). \quad (5.31)$$

For points  $\mathbf{w}$  that are sufficiently close to  $\hat{\mathbf{w}}$ , these expressions will give reasonable approximations for the error and its gradient.

Consider the particular case of a local quadratic approximation around a point  $\mathbf{w}^*$  that is a minimum of the error function. In this case there is no linear term, because  $\nabla E = 0$  at  $\mathbf{w}^*$ , and (5.28) becomes

$$E(\mathbf{w}) \underset{\approx}{=} E(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \quad (5.32)$$

where the Hessian  $\mathbf{H}$  is evaluated at  $\mathbf{w}^*$ . In order to interpret this geometrically, consider the eigenvalue equation for the Hessian matrix

$$\mathbf{H}\mathbf{u}_i = \lambda_i \mathbf{u}_i \quad (5.33)$$

where the eigenvectors  $\mathbf{u}_i$  form a complete orthonormal set (Appendix C) so that

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}. \quad (5.34)$$

We now expand  $(\mathbf{w} - \mathbf{w}^*)$  as a linear combination of the eigenvectors in the form

$$\mathbf{w} - \mathbf{w}^* = \sum_i \alpha_i \mathbf{u}_i. \quad (5.35)$$

This can be regarded as a transformation of the coordinate system in which the origin is translated to the point  $\mathbf{w}^*$ , and the axes are rotated to align with the eigenvectors (through the orthogonal matrix whose columns are the  $\mathbf{u}_i$ ), and is discussed in more detail in Appendix C. Substituting (5.35) into (5.32), and using (5.33) and (5.34), allows the error function to be written in the form

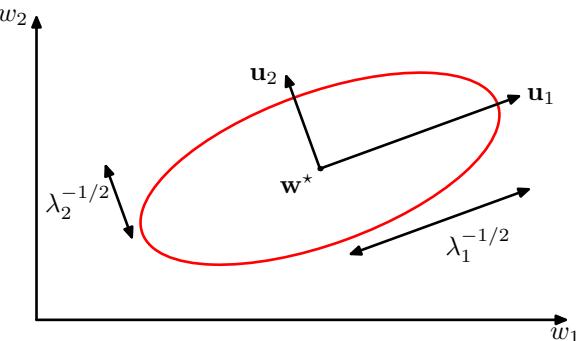
$$E(\mathbf{w}) = E(\mathbf{w}^*) + \frac{1}{2} \sum_i \lambda_i \alpha_i^2. \quad (5.36)$$

A matrix  $\mathbf{H}$  is said to be *positive definite* if, and only if,

$$\mathbf{v}^T \mathbf{H} \mathbf{v} > 0 \quad \text{for all } \mathbf{v} \neq \mathbf{0}. \quad (5.37)$$

**Figure 5.6** In the neighbourhood of a minimum  $\mathbf{w}^*$ , the error function can be approximated by a quadratic. Contours of constant error are then ellipses whose axes are aligned with the eigenvectors  $\mathbf{u}_i$  of the Hessian matrix, with lengths that are inversely proportional to the square roots of the corresponding eigenvalues  $\lambda_i$ .

eigenvalues



Because the eigenvectors  $\{\mathbf{u}_i\}$  form a complete set, an arbitrary vector  $\mathbf{v}$  can be written in the form

$$\mathbf{v} = \sum_i c_i \mathbf{u}_i. \quad (5.38)$$

From (5.33) and (5.34), we then have

$$\mathbf{v}^T \mathbf{H} \mathbf{v} = \sum_i c_i^2 \lambda_i \quad (5.39)$$

strictly

**Exercise 5.10**

**Exercise 5.11**

**Exercise 5.12**

and so  $\mathbf{H}$  will be positive definite if, and only if, all of its eigenvalues are positive. In the new coordinate system, whose basis vectors are given by the eigenvectors  $\{\mathbf{u}_i\}$ , the contours of constant  $E$  are ellipses centred on the origin, as illustrated in Figure 5.6. For a one-dimensional weight space, a stationary point  $w^*$  will be a minimum if

$$\left. \frac{\partial^2 E}{\partial w^2} \right|_{w^*} > 0. \quad (5.40)$$

The corresponding result in  $D$ -dimensions is that the Hessian matrix, evaluated at  $\mathbf{w}^*$ , should be positive definite.

### 5.2.3 Use of gradient information

As we shall see in Section 5.3, it is possible to evaluate the gradient of an error function efficiently by means of the backpropagation procedure. The use of this gradient information can lead to significant improvements in the speed with which the minima of the error function can be located. We can see why this is so, as follows.

In the quadratic approximation to the error function, given in (5.28), the error surface is specified by the quantities  $\mathbf{b}$  and  $\mathbf{H}$ , which contain a total of  $W(W + 3)/2$  independent elements (because the matrix  $\mathbf{H}$  is symmetric), where  $W$  is the dimensionality of  $\mathbf{w}$  (i.e., the total number of adaptive parameters in the network). The location of the minimum of this quadratic approximation therefore depends on  $O(W^2)$  parameters, and we should not expect to be able to locate the minimum until we have gathered  $O(W^2)$  independent pieces of information. If we do not make use of gradient information, we would expect to have to perform  $O(W^2)$  function

**Exercise 5.13**

evaluations, each of which would require  $O(W)$  steps. Thus, the computational effort needed to find the minimum using such an approach would be  $O(W^3)$ .

Now compare this with an algorithm that makes use of the gradient information. Because each evaluation of  $\nabla E$  brings  $W$  items of information, we might hope to find the minimum of the function in  $O(W)$  gradient evaluations. As we shall see, by using error backpropagation, each such evaluation takes only  $O(W)$  steps and so the minimum can now be found in  $O(W^2)$  steps. For this reason, the use of gradient information forms the basis of practical algorithms for training neural networks.

### 5.2.4 Gradient descent optimization

The simplest approach to using gradient information is to choose the weight update in (5.27) to comprise a small step in the direction of the negative gradient, so that

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)}) \quad (5.41)$$

where the parameter  $\eta > 0$  is known as the *learning rate*. After each such update, the gradient is re-evaluated for the new weight vector and the process repeated. Note that the error function is defined with respect to a training set, and so each step requires that the entire training set be processed in order to evaluate  $\nabla E$ . Techniques that use the whole data set at once are called *batch methods*. At each step the weight vector is moved in the direction of the greatest rate of decrease of the error function, and so this approach is known as *gradient descent or steepest descent*. Although such an approach might intuitively seem reasonable, in fact it turns out to be a poor algorithm, for reasons discussed in Bishop and Nabney (2008).

For batch optimization, there are more efficient methods, such as *conjugate gradients* and *quasi-Newton methods*, which are much more robust and much faster than simple gradient descent (Gill *et al.*, 1981; Fletcher, 1987; Nocedal and Wright, 1999). Unlike gradient descent, these algorithms have the property that the error function always decreases at each iteration unless the weight vector has arrived at a local or global minimum.

In order to find a sufficiently good minimum, it may be necessary to run a gradient-based algorithm multiple times, each time using a different randomly chosen starting point, and comparing the resulting performance on an independent validation set.

There is, however, an on-line version of gradient descent that has proved useful in practice for training neural networks on large data sets (Le Cun *et al.*, 1989). Error functions based on maximum likelihood for a set of independent observations comprise a sum of terms, one for each data point

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}). \quad (5.42)$$

On-line gradient descent, also known as *sequential gradient descent* or *stochastic gradient descent*, makes an update to the weight vector based on one data point at a time, so that

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)}). \quad (5.43)$$

This update is repeated by cycling through the data either in sequence or by selecting points at random with replacement. There are of course intermediate scenarios in which the updates are based on batches of data points.

One advantage of on-line methods compared to batch methods is that the former handle redundancy in the data much more efficiently. To see this, consider an extreme example in which we take a data set and double its size by duplicating every data point. Note that this simply multiplies the error function by a factor of 2 and so is equivalent to using the original error function. Batch methods will require double the computational effort to evaluate the batch error function gradient, whereas on-line methods will be unaffected. Another property of on-line gradient descent is the possibility of escaping from local minima, since a stationary point with respect to the error function for the whole data set will generally not be a stationary point for each data point individually.

Nonlinear optimization algorithms, and their practical application to neural network training, are discussed in detail in Bishop and Nabney (2008).

### 5.3. Error Backpropagation

---

Our goal in this section is to find an efficient technique for evaluating the gradient of an error function  $E(\mathbf{w})$  for a feed-forward neural network. We shall see that this can be achieved using a local message passing scheme in which information is sent alternately forwards and backwards through the network and is known as *error backpropagation*, or sometimes simply as *backprop*.

It should be noted that the term backpropagation is used in the neural computing literature to mean a variety of different things. For instance, the multilayer perceptron architecture is sometimes called a backpropagation network. The term backpropagation is also used to describe the training of a multilayer perceptron using gradient descent applied to a sum-of-squares error function. In order to clarify the terminology, it is useful to consider the nature of the training process more carefully. Most training algorithms involve an iterative procedure for minimization of an error function, with adjustments to the weights being made in a sequence of steps. At each such step, we can distinguish between two distinct stages. In the first stage, the derivatives of the error function with respect to the weights must be evaluated. As we shall see, the important contribution of the backpropagation technique is in providing a computationally efficient method for evaluating such derivatives. Because it is at this stage that errors are propagated backwards through the network, we shall use the term backpropagation specifically to describe the evaluation of derivatives. In the second stage, the derivatives are then used to compute the adjustments to be made to the weights. The simplest such technique, and the one originally considered by Rumelhart *et al.* (1986), involves gradient descent. It is important to recognize that the two stages are distinct. Thus, the first stage, namely the propagation of errors backwards through the network in order to evaluate derivatives, can be applied to many other kinds of network and not just the multilayer perceptron. It can also be applied to error functions other than just the simple sum-of-squares, and to the eval-

uation of other derivatives such as the Jacobian and Hessian matrices, as we shall see later in this chapter. Similarly, the second stage of weight adjustment using the calculated derivatives can be tackled using a variety of optimization schemes, many of which are substantially more powerful than simple gradient descent.

### 5.3.1 Evaluation of error-function derivatives

We now derive the backpropagation algorithm for a general network having arbitrary feed-forward topology, arbitrary differentiable nonlinear activation functions, and a broad class of error function. The resulting formulae will then be illustrated using a simple layered network structure having a single layer of sigmoidal hidden units together with a sum-of-squares error.

Many error functions of practical interest, for instance those defined by maximum likelihood for a set of i.i.d. data, comprise a sum of terms, one for each data point in the training set, so that

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}). \quad (5.44)$$

Here we shall consider the problem of evaluating  $\nabla E_n(\mathbf{w})$  for one such term in the error function. This may be used directly for sequential optimization, or the results can be accumulated over the training set in the case of batch methods.

Consider first a simple linear model in which the outputs  $y_k$  are linear combinations of the input variables  $x_i$  so that

$$y_k = \sum_i w_{ki} x_i \quad (5.45)$$

together with an error function that, for a particular input pattern  $n$ , takes the form

$$E_n = \frac{1}{2} \sum_k (y_{nk} - t_{nk})^2 \quad (5.46)$$

where  $y_{nk} = y_k(\mathbf{x}_n, \mathbf{w})$ . The gradient of this error function with respect to a weight  $w_{ji}$  is given by

$$\frac{\partial E_n}{\partial w_{ji}} = (y_{nj} - t_{nj}) x_{ni} \quad (5.47)$$

which can be interpreted as a ‘local’ computation involving the product of an ‘error signal’  $y_{nj} - t_{nj}$  associated with the output end of the link  $w_{ji}$  and the variable  $x_{ni}$  associated with the input end of the link. In Section 4.3.2, we saw how a similar formula arises with the logistic sigmoid activation function together with the cross entropy error function, and similarly for the softmax activation function together with its matching cross-entropy error function. We shall now see how this simple result extends to the more complex setting of multilayer feed-forward networks.

In a general feed-forward network, each unit computes a weighted sum of its inputs of the form

$$a_j = \sum_i w_{ji} z_i \quad (5.48)$$

where  $z_i$  is the activation of a unit, or input, that sends a connection to unit  $j$ , and  $w_{ji}$  is the weight associated with that connection. In Section 5.1, we saw that biases can be included in this sum by introducing an extra unit, or input, with activation fixed at +1. We therefore do not need to deal with biases explicitly. The sum in (5.48) is transformed by a nonlinear activation function  $h(\cdot)$  to give the activation  $z_j$  of unit  $j$  in the form

$$z_j = h(a_j). \quad (5.49)$$

Note that one or more of the variables  $z_i$  in the sum in (5.48) could be an input, and similarly, the unit  $j$  in (5.49) could be an output.

For each pattern in the training set, we shall suppose that we have supplied the corresponding input vector to the network and calculated the activations of all of the hidden and output units in the network by successive application of (5.48) and (5.49). This process is often called *forward propagation* because it can be regarded as a forward flow of information through the network.

Now consider the evaluation of the derivative of  $E_n$  with respect to a weight  $w_{ji}$ . The outputs of the various units will depend on the particular input pattern  $n$ . However, in order to keep the notation uncluttered, we shall omit the subscript  $n$  from the network variables. First we note that  $E_n$  depends on the weight  $w_{ji}$  only via the summed input  $a_j$  to unit  $j$ . We can therefore apply the chain rule for partial derivatives to give

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}}. \quad (5.50)$$

We now introduce a useful notation

$$\delta_j = \frac{\partial E_n}{\partial a_j} \quad (5.51)$$

where the  $\delta$ 's are often referred to as *errors* for reasons we shall see shortly. Using (5.48), we can write

$$\frac{\partial a_j}{\partial w_{ji}} = z_i. \quad (5.52)$$

Substituting (5.51) and (5.52) into (5.50), we then obtain

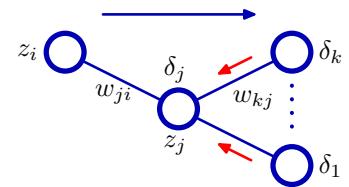
$$\frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i. \quad (5.53)$$

Equation (5.53) tells us that the required derivative is obtained simply by multiplying the value of  $\delta$  for the unit at the output end of the weight by the value of  $z$  for the unit at the input end of the weight (where  $z = 1$  in the case of a bias). Note that this takes the same form as for the simple linear model considered at the start of this section. Thus, in order to evaluate the derivatives, we need only to calculate the value of  $\delta_j$  for each hidden and output unit in the network, and then apply (5.53).

As we have seen already, for the output units, we have

$$\delta_k = y_k - t_k \quad (5.54)$$

**Figure 5.7** Illustration of the calculation of  $\delta_j$  for hidden unit  $j$  by backpropagation of the  $\delta$ 's from those units  $k$  to which unit  $j$  sends connections. The blue arrow denotes the direction of information flow during forward propagation, and the red arrows indicate the backward propagation of error information.



provided we are using the canonical link as the output-unit activation function. To evaluate the  $\delta$ 's for hidden units, we again make use of the chain rule for partial derivatives,

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} \quad (5.55)$$

where the sum runs over all units  $k$  to which unit  $j$  sends connections. The arrangement of units and weights is illustrated in Figure 5.7. Note that the units labelled  $k$  could include other hidden units and/or output units. In writing down (5.55), we are making use of the fact that variations in  $a_j$  give rise to variations in the error function only through variations in the variables  $a_k$ . If we now substitute the definition of  $\delta$  given by (5.51) into (5.55), and make use of (5.48) and (5.49), we obtain the following *backpropagation* formula

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k \quad (5.56)$$

which tells us that the value of  $\delta$  for a particular hidden unit can be obtained by propagating the  $\delta$ 's backwards from units higher up in the network, as illustrated in Figure 5.7. Note that the summation in (5.56) is taken over the first index on  $w_{kj}$  (corresponding to backward propagation of information through the network), whereas in the forward propagation equation (5.10) it is taken over the second index. Because we already know the values of the  $\delta$ 's for the output units, it follows that by recursively applying (5.56) we can evaluate the  $\delta$ 's for all of the hidden units in a feed-forward network, regardless of its topology.

The backpropagation procedure can therefore be summarized as follows.

### Error Backpropagation

1. Apply an input vector  $\mathbf{x}_n$  to the network and forward propagate through the network using (5.48) and (5.49) to find the activations of all the hidden and output units.
2. Evaluate the  $\delta_k$  for all the output units using (5.54).
3. Backpropagate the  $\delta$ 's using (5.56) to obtain  $\delta_j$  for each hidden unit in the network.
4. Use (5.53) to evaluate the required derivatives.

For batch methods, the derivative of the total error  $E$  can then be obtained by repeating the above steps for each pattern in the training set and then summing over all patterns:

$$\frac{\partial E}{\partial w_{ji}} = \sum_n \frac{\partial E_n}{\partial w_{ji}}. \quad (5.57)$$

In the above derivation we have implicitly assumed that each hidden or output unit in the network has the same activation function  $h(\cdot)$ . The derivation is easily generalized, however, to allow different units to have individual activation functions, simply by keeping track of which form of  $h(\cdot)$  goes with which unit.

### 5.3.2 A simple example

The above derivation of the backpropagation procedure allowed for general forms for the error function, the activation functions, and the network topology. In order to illustrate the application of this algorithm, we shall consider a particular example. This is chosen both for its simplicity and for its practical importance, because many applications of neural networks reported in the literature make use of this type of network. Specifically, we shall consider a two-layer network of the form illustrated in Figure 5.1, together with a sum-of-squares error, in which the output units have linear activation functions, so that  $y_k = a_k$ , while the hidden units have logistic sigmoid activation functions given by

sigmoidal

$$h(a) \equiv \tanh(a) \quad (5.58)$$

where

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}. \quad (5.59)$$

A useful feature of this function is that its derivative can be expressed in a particularly simple form:

$$h'(a) = 1 - h(a)^2. \quad (5.60)$$

We also consider a standard sum-of-squares error function, so that for pattern  $n$  the error is given by

$$E_n = \frac{1}{2} \sum_{k=1}^K (y_k - t_k)^2 \quad (5.61)$$

where  $y_k$  is the activation of output unit  $k$ , and  $t_k$  is the corresponding target, for a particular input pattern  $\mathbf{x}_n$ .

For each pattern in the training set in turn, we first perform a forward propagation using

$$a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i \quad (5.62)$$

$$z_j = \tanh(a_j) \quad (5.63)$$

$$y_k = \sum_{j=0}^M w_{kj}^{(2)} z_j. \quad (5.64)$$

Next we compute the  $\delta$ 's for each output unit using

$$\delta_k = y_k - t_k. \quad (5.65)$$

Then we backpropagate these to obtain  $\delta$ s for the hidden units using

$$\delta_j = (1 - z_j^2) \sum_{k=1}^K w_{kj} \delta_k. \quad (5.66)$$

Finally, the derivatives with respect to the first-layer and second-layer weights are given by

$$\frac{\partial E_n}{\partial w_{ji}^{(1)}} = \delta_j x_i, \quad \frac{\partial E_n}{\partial w_{kj}^{(2)}} = \delta_k z_j. \quad (5.67)$$

### 5.3.3 Efficiency of backpropagation

One of the most important aspects of backpropagation is its computational efficiency. To understand this, let us examine how the number of computer operations required to evaluate the derivatives of the error function scales with the total number  $W$  of weights and biases in the network. A single evaluation of the error function (for a given input pattern) would require  $O(W)$  operations, for sufficiently large  $W$ . This follows from the fact that, except for a network with very sparse connections, the number of weights is typically much greater than the number of units, and so the bulk of the computational effort in forward propagation is concerned with evaluating the sums in (5.48), with the evaluation of the activation functions representing a small overhead. Each term in the sum in (5.48) requires one multiplication and one addition, leading to an overall computational cost that is  $O(W)$ .

An alternative approach to backpropagation for computing the derivatives of the error function is to use finite differences. This can be done by perturbing each weight in turn, and approximating the derivatives by the expression

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji})}{\epsilon} + O(\epsilon) \quad (5.68)$$

where  $\epsilon \ll 1$ . In a software simulation, the accuracy of the approximation to the derivatives can be improved by making  $\epsilon$  smaller, until numerical roundoff problems arise. The accuracy of the finite differences method can be improved significantly by using symmetrical central differences of the form

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji} - \epsilon)}{2\epsilon} + O(\epsilon^2). \quad (5.69)$$

#### Exercise 5.14

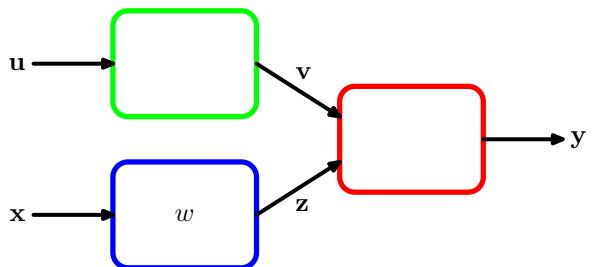
In this case, the  $O(\epsilon)$  corrections cancel, as can be verified by Taylor expansion on the right-hand side of (5.69), and so the residual corrections are  $O(\epsilon^2)$ . The number of computational steps is, however, roughly doubled compared with (5.68).

The main problem with numerical differentiation is that the highly desirable  $O(W)$  scaling has been lost. Each forward propagation requires  $O(W)$  steps, and

反向传播类似动态规划  
算法：迭代+存储共享  
用中间计算结果。

**Figure 5.8**

Illustration of a modular pattern recognition system in which the Jacobian matrix can be used to backpropagate error signals from the outputs through to earlier modules in the system.



there are  $W$  weights in the network each of which must be perturbed individually, so that the overall scaling is  $O(W^2)$ .

在实际操作中，有限差分计算的梯度结果可用于检查反向传播是否被正确实现以及其计算结果的准确性。

However, numerical differentiation plays an important role in practice, because a comparison of the derivatives calculated by backpropagation with those obtained using central differences provides a powerful check on the correctness of any software implementation of the backpropagation algorithm. When training networks in practice, derivatives should be evaluated using backpropagation, because this gives the greatest accuracy and numerical efficiency. However, the results should be compared with numerical differentiation using (5.69) for some test cases in order to check the correctness of the implementation.

### 5.3.4 The Jacobian matrix

We have seen how the derivatives of an error function with respect to the weights can be obtained by the propagation of errors backwards through the network. The technique of backpropagation can also be applied to the calculation of other derivatives. Here we consider the evaluation of the *Jacobian matrix*, whose elements are given by the derivatives of the network outputs with respect to the inputs

$$J_{ki} \equiv \frac{\partial y_k}{\partial x_i} \quad (5.70)$$

where each such derivative is evaluated with all other inputs held fixed. Jacobian matrices play a useful role in systems built from a number of distinct modules, as illustrated in Figure 5.8. Each module can comprise a fixed or adaptive function, which can be linear or nonlinear, so long as it is differentiable. Suppose we wish to minimize an error function  $E$  with respect to the parameter  $w$  in Figure 5.8. The derivative of the error function is given by

$$\frac{\partial E}{\partial w} = \sum_{k,j} \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_j} \frac{\partial z_j}{\partial w} \quad (5.71)$$

in which the Jacobian matrix for the red module in Figure 5.8 appears in the middle term.

Because the Jacobian matrix provides a measure of the local sensitivity of the outputs to changes in each of the input variables, it also allows any known errors  $\Delta x_i$

associated with the inputs to be propagated through the trained network in order to estimate their contribution  $\Delta y_k$  to the errors at the outputs, through the relation

$$\Delta y_k \simeq \sum_i \frac{\partial y_k}{\partial x_i} \Delta x_i \quad (5.72)$$

which is valid provided the  $|\Delta x_i|$  are small. In general, the network mapping represented by a trained neural network will be nonlinear, and so the elements of the Jacobian matrix will not be constants but will depend on the particular input vector used. Thus (5.72) is valid only for small perturbations of the inputs, and the Jacobian itself must be re-evaluated for each new input vector.

The Jacobian matrix can be evaluated using a backpropagation procedure that is similar to the one derived earlier for evaluating the derivatives of an error function with respect to the weights. We start by writing the element  $J_{ki}$  in the form

$$\begin{aligned} J_{ki} = \frac{\partial y_k}{\partial x_i} &= \sum_j \frac{\partial y_k}{\partial a_j} \frac{\partial a_j}{\partial x_i} \\ &= \sum_j w_{ji} \frac{\partial y_k}{\partial a_j} \end{aligned} \quad (5.73)$$

where we have made use of (5.48). The sum in (5.73) runs over all units  $j$  to which the input unit  $i$  sends connections (for example, over all units in the first hidden layer in the layered topology considered earlier). We now write down a recursive backpropagation formula to determine the derivatives  $\partial y_k / \partial a_j$

$$\begin{aligned} \frac{\partial y_k}{\partial a_j} &= \sum_l \frac{\partial y_k}{\partial a_l} \frac{\partial a_l}{\partial a_j} \\ &= h'(a_j) \sum_l w_{lj} \frac{\partial y_k}{\partial a_l} \end{aligned} \quad (5.74)$$

where the sum runs over all units  $l$  to which unit  $j$  sends connections (corresponding to the first index of  $w_{lj}$ ). Again, we have made use of (5.48) and (5.49). This backpropagation starts at the output units for which the required derivatives can be found directly from the functional form of the output-unit activation function. For instance, if we have individual sigmoidal activation functions at each output unit, then

$$\frac{\partial y_k}{\partial a_j} = \delta_{kj} \sigma'(a_j) \quad (5.75)$$

whereas for softmax outputs we have

$$\frac{\partial y_k}{\partial a_j} = \delta_{kj} y_k - y_k y_{kj} \quad (5.76)$$

We can summarize the procedure for evaluating the Jacobian matrix as follows. Apply the input vector corresponding to the point in input space at which the Jacobian matrix is to be found, and forward propagate in the usual way to obtain the

activations of all of the hidden and output units in the network. Next, for each row  $k$  of the Jacobian matrix, corresponding to the output unit  $k$ , backpropagate using the recursive relation (5.74), starting with (5.75) or (5.76), for all of the hidden units in the network. Finally, use (5.73) to do the backpropagation to the inputs. The Jacobian can also be evaluated using an alternative *forward* propagation formalism, which can be derived in an analogous way to the backpropagation approach given here.

Again, the implementation of such algorithms can be checked by using numerical differentiation in the form

$$\frac{\partial y_k}{\partial x_i} = \frac{y_k(x_i + \epsilon) - y_k(x_i - \epsilon)}{2\epsilon} + O(\epsilon^2) \quad (5.77)$$

which involves  $2D$  forward propagations for a network having  $D$  inputs.

## 5.4. The Hessian Matrix

We have shown how the technique of backpropagation can be used to obtain the first derivatives of an error function with respect to the weights in the network. Backpropagation can also be used to evaluate the second derivatives of the error, given by

$$\frac{\partial^2 E}{\partial w_{ji} \partial w_{lk}}. \quad (5.78)$$

Note that it is sometimes convenient to consider all of the weight and bias parameters as elements  $w_i$  of a single vector, denoted  $\mathbf{w}$ , in which case the second derivatives form the elements  $H_{ij}$  of the *Hessian* matrix  $\mathbf{H}$ , where  $i, j \in \{1, \dots, W\}$  and  $W$  is the total number of weights and biases. The Hessian plays an important role in many aspects of neural computing, including the following:

1. Several nonlinear optimization algorithms used for training neural networks are based on considerations of the second-order properties of the error surface, which are controlled by the Hessian matrix (Bishop and Nabney, 2008).
2. The Hessian forms the basis of a fast procedure for re-training a feed-forward network following a small change in the training data (Bishop, 1991).
3. The inverse of the Hessian has been used to identify the least significant weights in a network as part of network ‘pruning’ algorithms (Le Cun *et al.*, 1990).
4. The Hessian plays a central role in the Laplace approximation for a Bayesian neural network (see Section 5.7). Its inverse is used to determine the predictive distribution for a trained network, its eigenvalues determine the values of hyperparameters, and its determinant is used to evaluate the model evidence.

Various approximation schemes have been used to evaluate the Hessian matrix for a neural network. However, the Hessian can also be calculated exactly using an extension of the backpropagation technique.

在参数W的梯度时也可  
用类似的方法计算，但

Exercise 5.15

梯度是逐式法则，只  
不过相对反向传播，前向  
计算麻烦且效率低。

An important consideration for many applications of the Hessian is the efficiency with which it can be evaluated. If there are  $W$  parameters (weights and biases) in the network, then the Hessian matrix has dimensions  $W \times W$  and so the computational effort needed to evaluate the Hessian will scale like  $O(W^2)$  for each pattern in the data set. As we shall see, there are efficient methods for evaluating the Hessian whose scaling is indeed  $O(W^2)$ .

### 5.4.1 Diagonal approximation

Some of the applications for the Hessian matrix discussed above require the inverse of the Hessian, rather than the Hessian itself. For this reason, there has been some interest in using a diagonal approximation to the Hessian, in other words one that simply replaces the off-diagonal elements with zeros, because its inverse is trivial to evaluate. Again, we shall consider an error function that consists of a sum of terms, one for each pattern in the data set, so that  $E = \sum_n E_n$ . The Hessian can then be obtained by considering one pattern at a time, and then summing the results over all patterns. From (5.48), the diagonal elements of the Hessian, for pattern  $n$ , can be written

$$\frac{\partial^2 E_n}{\partial w_{ji}^2} = \frac{\partial^2 E_n}{\partial a_j^2} z_i^2. \quad (5.79)$$

Using (5.48) and (5.49), the second derivatives on the right-hand side of (5.79) can be found recursively using the chain rule of differential calculus to give a backpropagation equation of the form

$$\frac{\partial^2 E_n}{\partial a_j^2} = h'(a_j)^2 \sum_k \sum_{k'} w_{kj} w_{k'j} \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}} + h''(a_j) \sum_k w_{kj} \frac{\partial E_n}{\partial a_k}. \quad (5.80)$$

If we now neglect off-diagonal elements in the second-derivative terms, we obtain (Becker and Le Cun, 1989; Le Cun *et al.*, 1990)

$$\frac{\partial^2 E_n}{\partial a_j^2} = h'(a_j)^2 \sum_k w_{kj}^2 \frac{\partial^2 E_n}{\partial a_k^2} + h''(a_j) \sum_k w_{kj} \frac{\partial E_n}{\partial a_k}. \quad (5.81)$$

Note that the number of computational steps required to evaluate this approximation is  $O(W)$ , where  $W$  is the total number of weight and bias parameters in the network, compared with  $O(W^2)$  for the full Hessian.

Ricotti *et al.* (1988) also used the diagonal approximation to the Hessian, but they retained all terms in the evaluation of  $\partial^2 E_n / \partial a_j^2$  and so obtained exact expressions for the diagonal terms. Note that this no longer has  $O(W)$  scaling. The major problem with diagonal approximations, however, is that in practice the Hessian is typically found to be strongly nondiagonal, and so these approximations, which are driven mainly by computational convenience, must be treated with care.

by

### 5.4.2 Outer product approximation

[ When neural networks are applied to regression problems, it is common to use a sum-of-squares error function of the form ]

$$E = \frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2 \quad (5.82)$$

where we have considered the case of a single output in order to keep the notation simple (the extension to several outputs is straightforward). We can then write the Hessian matrix in the form

$$\mathbf{H} = \nabla \nabla E = \sum_{n=1}^N \nabla y_n (\nabla y_n)^T + \sum_{n=1}^N (y_n - t_n) \nabla \nabla y_n. \quad (5.83)$$

If the network has been trained on the data set, and its outputs  $y_n$  happen to be very close to the target values  $t_n$ , then the second term in (5.83) will be small and can be neglected. More generally, however, it may be appropriate to neglect this term by the following argument. Recall from Section 1.5.5 that the optimal function that minimizes a sum-of-squares loss is the conditional average of the target data. The quantity  $(y_n - t_n)$  is then a random variable with zero mean. If we assume that its value is uncorrelated with the value of the second derivative term on the right-hand side of (5.83), then the whole term will average to zero in the summation over  $n$ .

By neglecting the second term in (5.83), we arrive at the Levenberg–Marquardt approximation or outer product approximation (because the Hessian matrix is built up from a sum of outer products of vectors), given by

$$\mathbf{H} \simeq \sum_{n=1}^N \mathbf{b}_n \mathbf{b}_n^T \quad (5.84)$$

$\mathbf{b}_n \equiv \nabla a_n = \nabla y_n$  because the activation function for the output units is simply the identity. Evaluation of the outer product approximation for the Hessian is straightforward as it only involves first derivatives of the error function, which can be evaluated efficiently in  $O(W)$  steps using standard backpropagation. The elements of the matrix can then be found in  $O(W^2)$  steps by simple multiplication. It is important to emphasize that this approximation is only likely to be valid for a network that has been trained appropriately, and that for a general network mapping the second derivative terms on the right-hand side of (5.83) will typically not be negligible.]

[ In the case of the cross-entropy error function for a network with logistic sigmoid output-unit activation functions, the corresponding approximation is given by ]

$$\mathbf{H} \simeq \sum_{n=1}^N y_n(1-y_n) \mathbf{b}_n \mathbf{b}_n^T. \quad (5.85)$$

An analogous result can be obtained for multiclass networks having softmax output-unit activation functions.]

*Exercise 5.16*

*Exercise 5.17*

*Exercise 5.19*

*Exercise 5.20*

### 5.4.3 Inverse Hessian

We can use the outer-product approximation to develop a computationally efficient procedure for approximating the inverse of the Hessian (Hassibi and Stork, 1993). First we write the outer-product approximation in matrix notation as

$$\mathbf{H}_N = \sum_{n=1}^N \mathbf{b}_n \mathbf{b}_n^T \quad (5.86)$$

where  $\mathbf{b}_n \equiv \nabla_{\mathbf{w}} a_n$  is the contribution to the gradient of the output unit activation arising from data point  $n$ . We now derive a sequential procedure for building up the Hessian by including data points one at a time. Suppose we have already obtained the inverse Hessian using the first  $L$  data points. By separating off the contribution from data point  $L + 1$ , we obtain

$$\mathbf{H}_{L+1} = \mathbf{H}_L + \mathbf{b}_{L+1} \mathbf{b}_{L+1}^T. \quad (5.87)$$

In order to evaluate the inverse of the Hessian, we now consider the matrix identity

$$(\mathbf{M} + \mathbf{v}\mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{(\mathbf{M}^{-1}\mathbf{v})(\mathbf{v}^T\mathbf{M}^{-1})}{1 + \mathbf{v}^T\mathbf{M}^{-1}\mathbf{v}} \quad (5.88)$$

where  $\mathbf{I}$  is the unit matrix, which is simply a special case of the Woodbury identity (C.7). If we now identify  $\mathbf{H}_L$  with  $\mathbf{M}$  and  $\mathbf{b}_{L+1}$  with  $\mathbf{v}$ , we obtain

$$\mathbf{H}_{L+1}^{-1} = \mathbf{H}_L^{-1} - \frac{\mathbf{H}_L^{-1}\mathbf{b}_{L+1}\mathbf{b}_{L+1}^T\mathbf{H}_L^{-1}}{1 + \mathbf{b}_{L+1}^T\mathbf{H}_L^{-1}\mathbf{b}_{L+1}}. \quad (5.89)$$

In this way, data points are sequentially absorbed until  $L + 1 = N$  and the whole data set has been processed. This result therefore represents a procedure for evaluating the inverse of the Hessian using a single pass through the data set. The initial matrix  $\mathbf{H}_0$  is chosen to be  $\alpha\mathbf{I}$ , where  $\alpha$  is a small quantity, so that the algorithm actually finds the inverse of  $\mathbf{H} + \alpha\mathbf{I}$ . The results are not particularly sensitive to the precise value of  $\alpha$ . Extension of this algorithm to networks having more than one output is straightforward.

*Exercise 5.21*

We note here that the Hessian matrix can sometimes be calculated indirectly as part of the network training algorithm. In particular, quasi-Newton nonlinear optimization algorithms gradually build up an approximation to the inverse of the Hessian during training. Such algorithms are discussed in detail in Bishop and Nabney (2008).

### 5.4.4 Finite differences

As in the case of the first derivatives of the error function, we can find the second derivatives by using finite differences, with accuracy limited by numerical precision. If we perturb each possible pair of weights in turn, we obtain

$$\begin{aligned} \frac{\partial^2 E}{\partial w_{ji} \partial w_{lk}} &= \frac{1}{4\epsilon^2} \{E(w_{ji} + \epsilon, w_{lk} + \epsilon) - E(w_{ji} + \epsilon, w_{lk} - \epsilon) \\ &\quad - E(w_{ji} - \epsilon, w_{lk} + \epsilon) + E(w_{ji} - \epsilon, w_{lk} - \epsilon)\} + O(\epsilon^2). \end{aligned} \quad (5.90)$$

課一：

Again, by using a symmetrical central differences formulation, we ensure that the residual errors are  $O(\epsilon^2)$  rather than  $O(\epsilon)$ . Because there are  $W^2$  elements in the Hessian matrix, and because the evaluation of each element requires four forward propagations each needing  $O(W)$  operations (per pattern), we see that this approach will require  $O(W^3)$  operations to evaluate the complete Hessian. It therefore has poor scaling properties, although in practice it is very useful as a check on the software implementation of backpropagation methods.]

**方案二：** [ A more efficient version of numerical differentiation can be found by applying central differences to the first derivatives of the error function, which are themselves calculated using backpropagation. This gives

$$\frac{\partial^2 E}{\partial w_{ji} \partial w_{lk}} = \frac{1}{2\epsilon} \left\{ \frac{\partial E}{\partial w_{ji}}(w_{lk} + \epsilon) - \frac{\partial E}{\partial w_{ji}}(w_{lk} - \epsilon) \right\} + O(\epsilon^2). \quad (5.91)$$

Because there are now only  $W$  weights to be perturbed, and because the gradients can be evaluated in  $O(W)$  steps, we see that this method gives the Hessian in  $O(W^2)$  operations.]

#### 5.4.5 Exact evaluation of the Hessian

So far, we have considered various approximation schemes for evaluating the Hessian matrix or its inverse. The Hessian can also be evaluated exactly, for a network of arbitrary feed-forward topology, using extension of the technique of back-propagation used to evaluate first derivatives, which shares many of its desirable features including computational efficiency (Bishop, 1991; Bishop, 1992). It can be applied to any differentiable error function that can be expressed as a function of the network outputs and to networks having arbitrary differentiable activation functions. The number of computational steps needed to evaluate the Hessian scales like  $O(W^2)$ . Similar algorithms have also been considered by Buntine and Weigend (1993).

Here we consider the specific case of a network having two layers of weights, for which the required equations are easily derived. We shall use indices  $i$  and  $i'$  to denote inputs, indices  $j$  and  $j'$  to denote hidden units, and indices  $k$  and  $k'$  to denote outputs. {We first define

$$\delta_k = \frac{\partial E_n}{\partial a_k}, \quad M_{kk'} \equiv \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}} \quad (5.92)$$

where  $E_n$  is the contribution to the error from data point  $n$ . The Hessian matrix for this network can then be considered in three separate blocks as follows.

1. Both weights in the second layer:

$$\frac{\partial^2 E_n}{\partial w_{kj}^{(2)} \partial w_{k'j'}^{(2)}} = z_j z_{j'} M_{kk'}. \quad (5.93)$$

2. Both weights in the first layer:

$$\begin{aligned} \frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{j'i'}^{(1)}} &= x_i x_{i'} h''(a_{j'}) I_{jj'} \sum_k w_{kj'}^{(2)} \delta_k \\ &+ x_i x_{i'} h'(a_{j'}) h'(a_j) \sum_k \sum_{k'} w_{k'j'}^{(2)} w_{kj}^{(2)} M_{kk'}. \end{aligned} \quad (5.94)$$

3. One weight in each layer:

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{kj'}^{(2)}} = x_i h'(a_j) \left\{ \delta_k I_{jk'} + z_j \sum_{k'} w_{k'j}^{(2)} \cancel{H_{kk'}} \right\}. \quad (5.95)$$

Here  $I_{jj'}$  is the  $j, j'$  element of the identity matrix. If one or both of the weights is a bias term, then the corresponding expressions are obtained simply by setting the appropriate activation(s) to 1. Inclusion of skip-layer connections is straightforward.

*Exercise 5.23*

#### 5.4.6 Fast multiplication by the Hessian

For many applications of the Hessian, the quantity of interest is not the Hessian matrix  $\mathbf{H}$  itself but the product of  $\mathbf{H}$  with some vector  $\mathbf{v}$ . We have seen that the evaluation of the Hessian takes  $O(W^2)$  operations, and it also requires storage that is  $O(W^2)$ . The vector  $\mathbf{v}^T \mathbf{H}$  that we wish to calculate, however, has only  $W$  elements, so instead of computing the Hessian as an intermediate step, we can instead try to find an efficient approach to evaluating  $\mathbf{v}^T \mathbf{H}$  directly in a way that requires only  $O(W)$  operations.

To do this, we first note that

$$\mathbf{v}^T \mathbf{H} = \mathbf{v}^T \nabla(\nabla E) \quad (5.96)$$

where  $\nabla$  denotes the gradient operator in weight space. We can then write down the standard forward-propagation and backpropagation equations for the evaluation of  $\nabla E$  and apply (5.96) to these equations to give a set of forward-propagation and backpropagation equations for the evaluation of  $\mathbf{v}^T \mathbf{H}$  (Møller, 1993; Pearlmutter, 1994). This corresponds to acting on the original forward-propagation and backpropagation equations with a differential operator  $\mathbf{v}^T \nabla$ . Pearlmutter (1994) used the notation  $\mathcal{R}\{\cdot\}$  to denote the operator  $\mathbf{v}^T \nabla$ , and we shall follow this convention. The analysis is straightforward and makes use of the usual rules of differential calculus, together with the result

$$\mathcal{R}\{\mathbf{w}\} = \mathbf{v}. \quad (5.97)$$

The technique is best illustrated with a simple example, and again we choose a two-layer network of the form shown in Figure 5.1, with linear output units and a sum-of-squares error function. As before, we consider the contribution to the error function from one pattern in the data set. The required vector is then obtained as

usual by summing over the contributions from each of the patterns separately. For the two-layer network, the forward-propagation equations are given by

$$a_j = \sum_i w_{ji}x_i \quad (5.98)$$

$$z_j = h(a_j) \quad (5.99)$$

$$y_k = \sum_j w_{kj}z_j. \quad (5.100)$$

① 前向传播：依次计算式(5.101)-(5.103)

We now act on these equations using the  $\mathcal{R}\{\cdot\}$  operator to obtain a set of forward propagation equations in the form

$$\mathcal{R}\{a_j\} = \sum_i v_{ji}x_i \quad (5.101)$$

$$\mathcal{R}\{z_j\} = h'(a_j)\mathcal{R}\{a_j\} \quad (5.102)$$

$$\mathcal{R}\{y_k\} = \sum_j w_{kj}\mathcal{R}\{z_j\} + \sum_j v_{kj}z_j \quad (5.103)$$

where  $v_{ji}$  is the element of the vector  $\mathbf{v}$  that corresponds to the weight  $w_{ji}$ . Quantities of the form  $\mathcal{R}\{z_j\}$ ,  $\mathcal{R}\{a_j\}$  and  $\mathcal{R}\{y_k\}$  are to be regarded as new variables whose values are found using the above equations.

Because we are considering a sum-of-squares error function, we have the following standard backpropagation expressions:

$$\delta_k = y_k - t_k \quad (5.104)$$

$$\delta_j = h'(a_j) \sum_k w_{kj}\delta_k. \quad (5.105)$$

② 后向传播，利用①中结果，依次计算式(5.106)-(5.107)

Again, we act on these equations with the  $\mathcal{R}\{\cdot\}$  operator to obtain a set of backpropagation equations in the form

-(5.107)

$$\mathcal{R}\{\delta_k\} = \mathcal{R}\{y_k\} \quad (5.106)$$

$$\begin{aligned} \mathcal{R}\{\delta_j\} &= h''(a_j)\mathcal{R}\{a_j\} \sum_k w_{kj}\delta_k \\ &\quad + h'(a_j) \sum_k v_{kj}\delta_k + h'(a_j) \sum_k w_{kj}\mathcal{R}\{\delta_k\}. \end{aligned} \quad (5.107)$$

③ 利用①、②所得结果，代入(5.110)-(5.113)中得最终结果

Finally, we have the usual equations for the first derivatives of the error

$$\frac{\partial E}{\partial w_{kj}} = \delta_k z_j \quad (5.108)$$

$$\frac{\partial E}{\partial w_{ji}} = \delta_j x_i \quad (5.109)$$

and acting on these with the  $\mathcal{R}\{\cdot\}$  operator, we obtain expressions for the elements of the vector  $\mathbf{v}^T \mathbf{H}$

$$\mathcal{R}\left\{\frac{\partial E}{\partial w_{kj}}\right\} = \mathcal{R}\{\delta_k\}z_j + \delta_k \mathcal{R}\{z_j\} \quad (5.110)$$

$$\mathcal{R}\left\{\frac{\partial E}{\partial w_{ji}}\right\} = x_i \mathcal{R}\{\delta_j\}. \quad (5.111)$$

The implementation of this algorithm involves the introduction of additional variables  $\mathcal{R}\{a_j\}$ ,  $\mathcal{R}\{z_j\}$  and  $\mathcal{R}\{\delta_j\}$  for the hidden units and  $\mathcal{R}\{\delta_k\}$  and  $\mathcal{R}\{y_k\}$  for the output units. For each input pattern, the values of these quantities can be found using the above results, and the elements of  $\mathbf{v}^T \mathbf{H}$  are then given by (5.110) and (5.111). An elegant aspect of this technique is that the equations for evaluating  $\mathbf{v}^T \mathbf{H}$  mirror closely those for standard forward and backward propagation, and so the extension of existing software to compute this product is typically straightforward.

If desired, the technique can be used to evaluate the full Hessian matrix by choosing the vector  $\mathbf{v}$  to be given successively by a series of unit vectors of the form  $(0, 0, \dots, 1, \dots, 0)$  each of which picks out one column of the Hessian. This leads to a formalism that is analytically equivalent to the backpropagation procedure of Bishop (1992), as described in Section 5.4.5, though with some loss of efficiency due to redundant calculations.

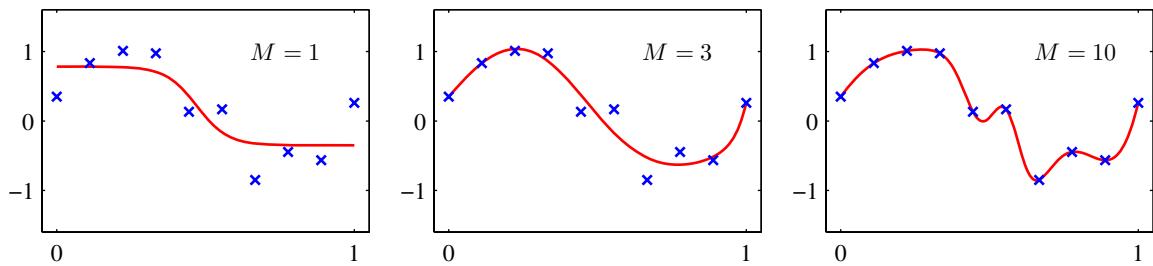
当V和类别one-hot  
编码的单位向量时，上  
述算法可用于精确计  
算H。只相比第  
5.4.5节中的算法，计算  
效率要低一点。

## 5.5. Regularization in Neural Networks

The number of input and output units in a neural network is generally determined by the dimensionality of the data set, whereas the number  $M$  of hidden units is a free parameter that can be adjusted to give the best predictive performance. Note that  $M$  controls the number of parameters (weights and biases) in the network, and so we might expect that in a maximum likelihood setting there will be an optimum value of  $M$  that gives the best generalization performance, corresponding to the optimum balance between under-fitting and over-fitting. Figure 5.9 shows an example of the effect of different values of  $M$  for the sinusoidal regression problem.

The generalization error, however, is not a simple function of  $M$  due to the presence of local minima in the error function, as illustrated in Figure 5.10. Here we see the effect of choosing multiple random initializations for the weight vector for a range of values of  $M$ . The overall best validation set performance in this case occurred for a particular solution having  $M = 8$ . In practice, one approach to choosing  $M$  is in fact to plot a graph of the kind shown in Figure 5.10 and then to choose the specific solution having the smallest validation set error.

There are, however, other ways to control the complexity of a neural network model in order to avoid over-fitting. From our discussion of polynomial curve fitting in Chapter 1, we see that an alternative approach is to choose a relatively large value for  $M$  and then to control complexity by the addition of a regularization term to the error function. The simplest regularizer is the quadratic, giving a regularized error



**Figure 5.9** Examples of two-layer networks trained on 10 data points drawn from the sinusoidal data set. The graphs show the result of fitting networks having  $M = 1, 3$  and 10 hidden units, respectively, by minimizing a sum-of-squares error function using a scaled conjugate-gradient algorithm.

of the form

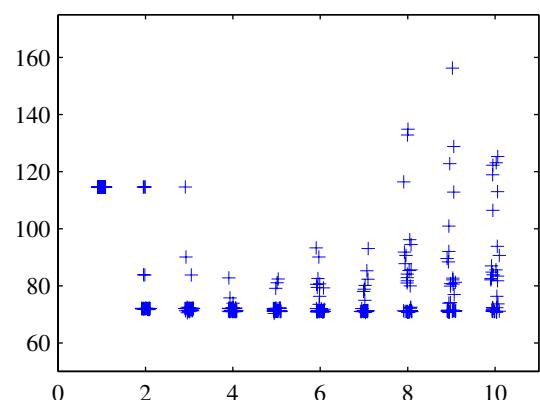
$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}. \quad (5.112)$$

This regularizer is also known as *weight decay* and has been discussed at length in Chapter 3. The effective model complexity is then determined by the choice of the regularization coefficient  $\lambda$ . As we have seen previously, this regularizer can be interpreted as the negative logarithm of a zero-mean Gaussian prior distribution over the weight vector  $\mathbf{w}$ .

### 5.5.1 Consistent Gaussian priors

One of the limitations of simple weight decay in the form (5.112) is that it is inconsistent with certain scaling properties of network mappings. To illustrate this, consider a multilayer perceptron network having two layers of weights and linear output units, which performs a mapping from a set of input variables  $\{x_i\}$  to a set of output variables  $\{y_k\}$ . The activations of the hidden units in the first hidden layer

**Figure 5.10** Plot of the sum-of-squares test-set error for the polynomial data set versus the number of hidden units in the network, with 30 random starts for each network size, showing the effect of local minima. For each new start, the weight vector was initialized by sampling from an isotropic Gaussian distribution having a mean of zero and a variance of 10.



take the form

$$z_j = h \left( \sum_i w_{ji} x_i + w_{j0} \right) \quad (5.113)$$

while the activations of the output units are given by

$$y_k = \sum_j w_{kj} z_j + w_{k0}. \quad (5.114)$$

Suppose we perform a linear transformation of the input data of the form

$$x_i \rightarrow \tilde{x}_i = ax_i + b. \quad (5.115)$$

Then we can arrange for the mapping performed by the network to be unchanged by making a corresponding linear transformation of the weights and biases from the inputs to the units in the hidden layer of the form

$$w_{ji} \rightarrow \tilde{w}_{ji} = \frac{1}{a} w_{ji} \quad (5.116)$$

$$w_{j0} \rightarrow \tilde{w}_{j0} = w_{j0} - \frac{b}{a} \sum_i w_{ji}. \quad (5.117)$$

Similarly, a linear transformation of the output variables of the network of the form

$$y_k \rightarrow \tilde{y}_k = cy_k + d \quad (5.118)$$

can be achieved by making a transformation of the second-layer weights and biases using

$$w_{kj} \rightarrow \tilde{w}_{kj} = cw_{kj} \quad (5.119)$$

$$w_{k0} \rightarrow \tilde{w}_{k0} = cw_{k0} + d. \quad (5.120)$$

If we train one network using the original data and one network using data for which the input and/or target variables are transformed by one of the above linear transformations, then consistency requires that we should obtain equivalent networks that differ only by the linear transformation of the weights as given. Any regularizer should be consistent with this property, otherwise it arbitrarily favours one solution over another, equivalent one. Clearly, simple weight decay (5.112), that treats all weights and biases on an equal footing, does not satisfy this property.

// We therefore look for a regularizer which is invariant under the linear transformations (5.116), (5.117), (5.119) and (5.120). These require that the regularizer should be invariant to re-scaling of the weights and to shifts of the biases. Such a regularizer is given by

$$\frac{\lambda_1}{2} \sum_{w \in \mathcal{W}_1} w^2 + \frac{\lambda_2}{2} \sum_{w \in \mathcal{W}_2} w^2 \quad (5.121)$$

where  $\mathcal{W}_1$  denotes the set of weights in the first layer,  $\mathcal{W}_2$  denotes the set of weights in the second layer, and biases are excluded from the summations. This regularizer

### Exercise 5.24

will remain unchanged under the weight transformations provided the regularization parameters are re-scaled using  $\lambda_1 \rightarrow a^{1/2}\lambda_1$  and  $\lambda_2 \rightarrow c^{-1/2}\lambda_2$ . 3

The regularizer (5.121) corresponds to a prior of the form

$$p(\mathbf{w}|\alpha_1, \alpha_2) \propto \exp\left(-\frac{\alpha_1}{2} \sum_{w \in \mathcal{W}_1} w^2 - \frac{\alpha_2}{2} \sum_{w \in \mathcal{W}_2} w^2\right). \quad (5.122)$$

包含偏置项  $w_0, w_k$ , 而  $\mathcal{W}_1, \mathcal{W}_2$  中不包含偏置

Note that priors of this form are *improper* (they cannot be normalized) because the bias parameters are unconstrained. 3 The use of improper priors can lead to difficulties in selecting regularization coefficients and in model comparison within the Bayesian framework, because the corresponding evidence is zero. It is therefore common to include separate priors for the biases (which then break shift invariance) having their own hyperparameters. 3 We can illustrate the effect of the resulting four hyperparameters by drawing samples from the prior and plotting the corresponding network functions, as shown in Figure 5.11.

More generally, we can consider priors in which the weights are divided into any number of groups  $\mathcal{W}_k$  so that

$$p(\mathbf{w}) \propto \exp\left(-\frac{1}{2} \sum_k \alpha_k \|\mathbf{w}\|_k^2\right) \quad (5.123)$$

where

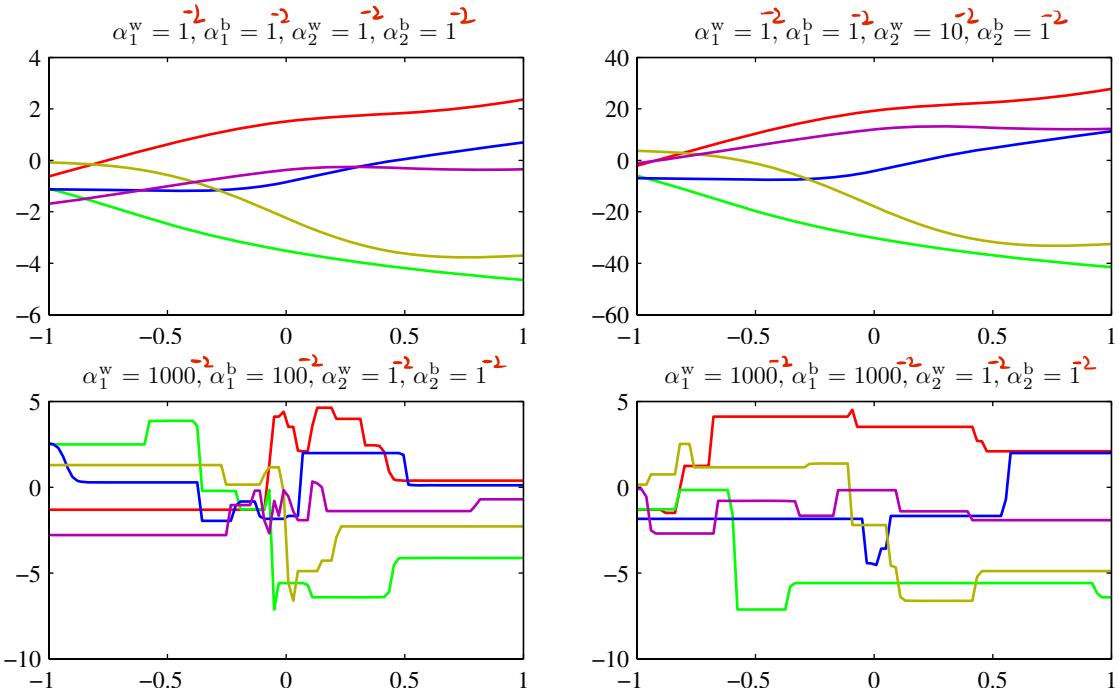
$$\|\mathbf{w}\|_k^2 = \sum_{j \in \mathcal{W}_k} w_j^2. \quad (5.124)$$

As a special case of this prior, if we choose the groups to correspond to the sets of weights associated with each of the input units, and we optimize the marginal likelihood with respect to the corresponding parameters  $\alpha_k$ , we obtain *automatic relevance determination* as discussed in Section 7.2.2. 3

## 5.5.2 Early stopping

An alternative to regularization as a way of controlling the effective complexity of a network is the procedure of *early stopping*. The training of nonlinear network models corresponds to an iterative reduction of the error function defined with respect to a set of training data. For many of the optimization algorithms used for network training, such as conjugate gradients, the error is a nonincreasing function of the iteration index. However, the error measured with respect to independent data, generally called a validation set, often shows a decrease at first, followed by an increase as the network starts to over-fit. Training can therefore be stopped at the point of smallest error with respect to the validation data set, as indicated in Figure 5.12, in order to obtain a network having good generalization performance.

The behaviour of the network in this case is sometimes explained qualitatively in terms of the effective number of degrees of freedom in the network, in which this number starts out small and then grows during the training process, corresponding to a steady increase in the effective complexity of the model. Halting training before

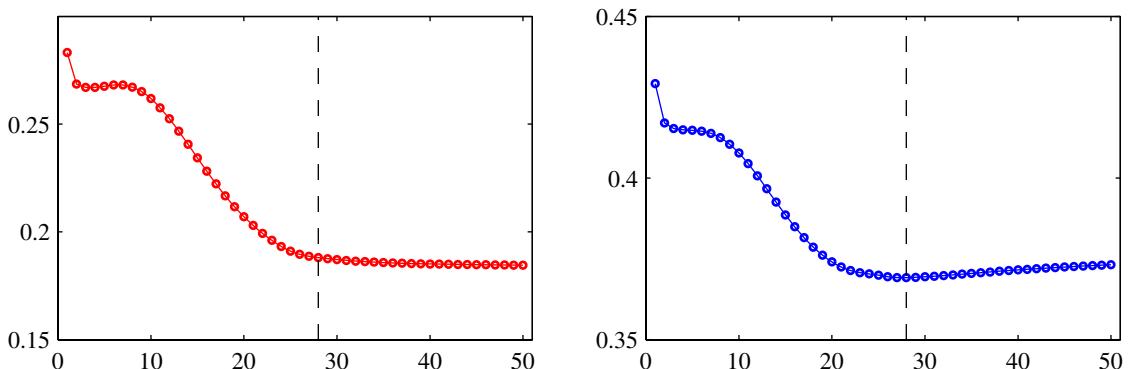


**Figure 5.11** Illustration of the effect of the hyperparameters governing the prior distribution over weights and biases in a two-layer network having a single input, a single linear output, and 12 hidden units having ‘tanh’ activation functions. The priors are governed by four hyperparameters  $\alpha_1^b$ ,  $\alpha_1^w$ ,  $\alpha_2^b$ , and  $\alpha_2^w$ , which represent the precisions of the Gaussian distributions of the first-layer biases, first-layer weights, second-layer biases, and second-layer weights, respectively. We see that the parameter  $\alpha_2^w$  governs the vertical scale of functions (note the different vertical axis ranges on the top two diagrams),  $\alpha_1^w$  governs the horizontal scale of variations in the function values, and  $\alpha_1^b$  governs the horizontal range over which variations occur. The parameter  $\alpha_2^b$ , whose effect is not illustrated here, governs the range of vertical offsets of the functions.

a minimum of the training error has been reached then represents a way of limiting the effective network complexity.

In the case of a quadratic error function, we can verify this insight, and show that early stopping should exhibit similar behaviour to regularization using a simple weight-decay term. This can be understood from Figure 5.13, in which the axes in weight space have been rotated to be parallel to the eigenvectors of the Hessian matrix. If, in the absence of weight decay, the weight vector starts at the origin and proceeds during training along a path that follows the local negative gradient vector, then the weight vector will move initially parallel to the  $w_2$  axis through a point corresponding roughly to  $\tilde{w}$  and then move towards the minimum of the error function  $w_{ML}$ . This follows from the shape of the error surface and the widely differing eigenvalues of the Hessian. Stopping at a point near  $\tilde{w}$  is therefore similar to weight decay. The relationship between early stopping and weight decay can be made quantitative, thereby showing that the quantity  $\tau\eta$  (where  $\tau$  is the iteration index, and  $\eta$  is the learning rate parameter) plays the role of the reciprocal of the regularization

### Exercise 5.25



**Figure 5.12** An illustration of the behaviour of training set error (left) and validation set error (right) during a typical training session, as a function of the iteration step, for the sinusoidal data set. The goal of achieving the best generalization performance suggests that training should be stopped at the point shown by the vertical dashed lines, corresponding to the minimum of the validation set error.

parameter  $\lambda$ . The effective number of parameters in the network therefore grows during the course of training.

### 5.5.3 Invariances

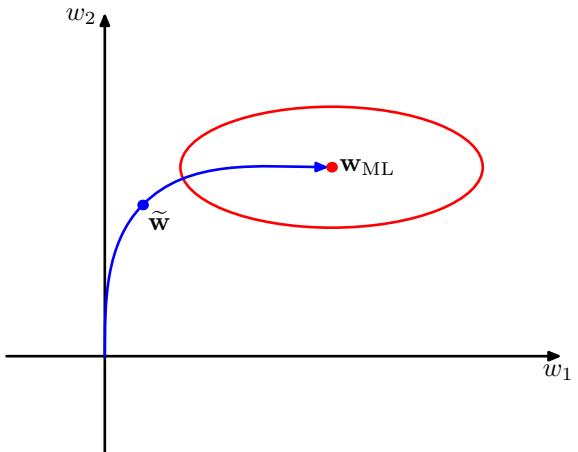
In many applications of pattern recognition, it is known that predictions should be unchanged, or *invariant*, under one or more transformations of the input variables. For example, in the classification of objects in two-dimensional images, such as handwritten digits, a particular object should be assigned the same classification irrespective of its position within the image (*translation invariance*) or of its size (*scale invariance*). Such transformations produce significant changes in the raw data, expressed in terms of the intensities at each of the pixels in the image, and yet should give rise to the same output from the classification system. Similarly in speech recognition, small levels of nonlinear warping along the time axis, which preserve temporal ordering, should not change the interpretation of the signal.

If sufficiently large numbers of training patterns are available, then an adaptive model such as a neural network can learn the invariance, at least approximately. This involves including within the training set a sufficiently large number of examples of the effects of the various transformations. Thus, for translation invariance in an image, the training set should include examples of objects at many different positions.

This approach may be impractical, however, if the number of training examples is limited, or if there are several invariants (because the number of combinations of transformations grows exponentially with the number of such transformations). We therefore seek alternative approaches for encouraging an adaptive model to exhibit the required invariances. These can broadly be divided into four categories:

1. The training set is augmented using replicas of the training patterns, transformed according to the desired invariances. For instance, in our digit recognition example, we could make multiple copies of each example in which the

**Figure 5.13** A schematic illustration of why early stopping can give similar results to weight decay in the case of a quadratic error function. The ellipse shows a contour of constant error, and  $w_{ML}$  denotes the minimum of the error function. If the weight vector starts at the origin and moves according to the local negative gradient direction, then it will follow the path shown by the curve. By stopping training early, a weight vector  $\tilde{w}$  is found that is qualitatively similar to that obtained with a simple weight-decay regularizer and training to the minimum of the regularized error, as can be seen by comparing with Figure 3.15.

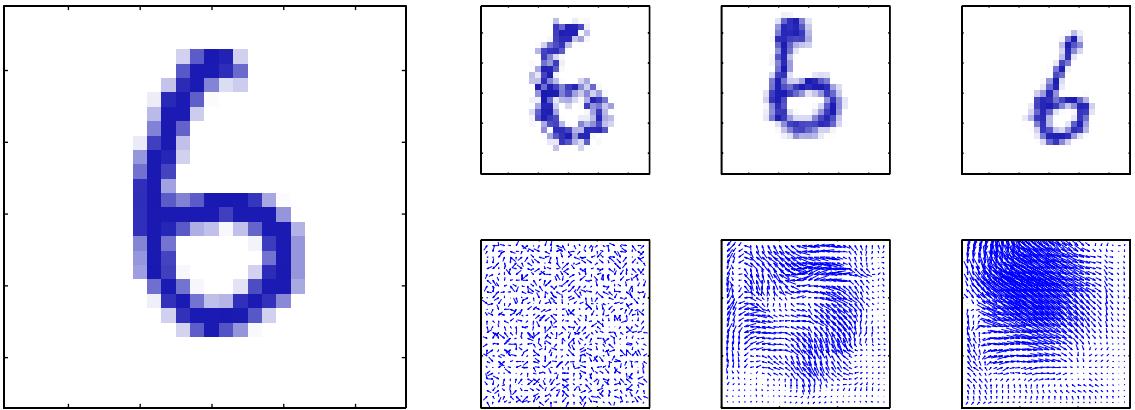


digit is shifted to a different position in each image.

2. A regularization term is added to the error function that penalizes changes in the model output when the input is transformed. This leads to the technique of *tangent propagation*, discussed in Section 5.5.4.
3. Invariance is built into the pre-processing by extracting features that are invariant under the required transformations. Any subsequent regression or classification system that uses such features as inputs will necessarily also respect these invariances.
4. The final option is to build the invariance properties into the structure of a neural network (or into the definition of a kernel function in the case of techniques such as the relevance vector machine). One way to achieve this is through the use of local receptive fields and shared weights, as discussed in the context of convolutional neural networks in Section 5.5.6.

Approach 1 is often relatively easy to implement and can be used to encourage complex invariances such as those illustrated in Figure 5.14. For sequential training algorithms, this can be done by transforming each input pattern before it is presented to the model so that, if the patterns are being recycled, a different transformation (drawn from an appropriate distribution) is added each time. For batch methods, a similar effect can be achieved by replicating each data point a number of times and transforming each copy independently. The use of such augmented data can lead to significant improvements in generalization (Simard *et al.*, 2003), although it can also be computationally costly.

Approach 2 leaves the data set unchanged but modifies the error function through the addition of a regularizer. In Section 5.5.5, we shall show that this approach is closely related to approach 1.



**Figure 5.14** Illustration of the synthetic warping of a handwritten digit. The original image is shown on the left. On the right, the top row shows three examples of warped digits, with the corresponding displacement fields shown on the bottom row. These displacement fields are generated by sampling random displacements  $\Delta x, \Delta y \in (0, 1)$  at each pixel and then smoothing by convolution with Gaussians of width 0.01, 30 and 60 respectively.

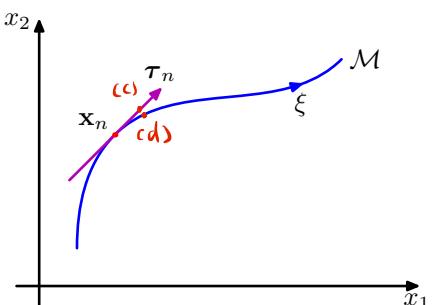
One advantage of approach 3 is that it can correctly extrapolate well beyond the range of transformations included in the training set. However, it can be difficult to find hand-crafted features with the required invariances that do not also discard information that can be useful for discrimination.

#### 5.5.4 Tangent propagation

We can use regularization to encourage models to be invariant to transformations of the input through the technique of *tangent propagation* (Simard *et al.*, 1992). Consider the effect of a transformation on a particular input vector  $x_n$ . Provided the transformation is continuous (such as translation or rotation, but not mirror reflection for instance), then the transformed pattern will sweep out a manifold  $\mathcal{M}$  within the  $D$ -dimensional input space. This is illustrated in Figure 5.15, for the case of  $D = 2$  for simplicity. Suppose the transformation is governed by a single parameter  $\xi$  (which might be rotation angle for instance). Then the subspace  $\mathcal{M}$  swept out by  $x_n$

tangent propagation  
切线传播

**Figure 5.15** Illustration of a two-dimensional input space showing the effect of a continuous transformation on a particular input vector  $x_n$ . A one-dimensional transformation, parameterized by the continuous variable  $\xi$ , applied to  $x_n$  causes it to sweep out a one-dimensional manifold  $\mathcal{M}$ . Locally, the effect of the transformation can be approximated by the tangent vector  $\tau_n$ .



will be one-dimensional, and will be parameterized by  $\xi$ . Let the vector that results from acting on  $\mathbf{x}_n$  by this transformation be denoted by  $\mathbf{s}(\mathbf{x}_n, \xi)$ , which is defined so that  $\mathbf{s}(\mathbf{x}, 0) = \mathbf{x}$ . Then the tangent to the curve  $\mathcal{M}$  is given by the directional derivative  $\boldsymbol{\tau} = \partial \mathbf{s} / \partial \xi$ , and the tangent vector at the point  $\mathbf{x}_n$  is given by

$$\boldsymbol{\tau}_n = \left. \frac{\partial \mathbf{s}(\mathbf{x}_n, \xi)}{\partial \xi} \right|_{\xi=0}. \quad (5.125)$$

Under a transformation of the input vector, the network output vector will, in general, change. The derivative of output  $k$  with respect to  $\xi$  is given by

$$\left. \frac{\partial y_k}{\partial \xi} \right|_{\xi=0} = \sum_{i=1}^D \left. \frac{\partial y_k}{\partial x_i} \frac{\partial x_i}{\partial \xi} \right|_{\xi=0} = \sum_{i=1}^D J_{ki} \tau_i \quad (5.126)$$

where  $J_{ki}$  is the  $(k, i)$  element of the Jacobian matrix  $\mathbf{J}$ , as discussed in Section 5.3.4. The result (5.126) can be used to modify the standard error function, so as to encourage local invariance in the neighbourhood of the data points, by the addition to the original error function  $E$  of a regularization function  $\Omega$  to give a total error function of the form

$$\tilde{E} = E + \lambda \Omega \quad (5.127)$$

where  $\lambda$  is a regularization coefficient and

$$\Omega = \frac{1}{2} \sum_n \sum_k \left( \left. \frac{\partial y_{nk}}{\partial \xi} \right|_{\xi=0} \right)^2 = \frac{1}{2} \sum_n \sum_k \left( \sum_{i=1}^D J_{nki} \tau_{ni} \right)^2. \quad (5.128)$$

The regularization function will be zero when the network mapping function is invariant under the transformation in the neighbourhood of each pattern vector, and the value of the parameter  $\lambda$  determines the balance between fitting the training data and learning the invariance property.

In a practical implementation, the tangent vector  $\boldsymbol{\tau}_n$  can be approximated using finite differences, by subtracting the original vector  $\mathbf{x}_n$  from the corresponding vector after transformation using a small value of  $\xi$ , and then dividing by  $\xi$ . This is illustrated in Figure 5.16.

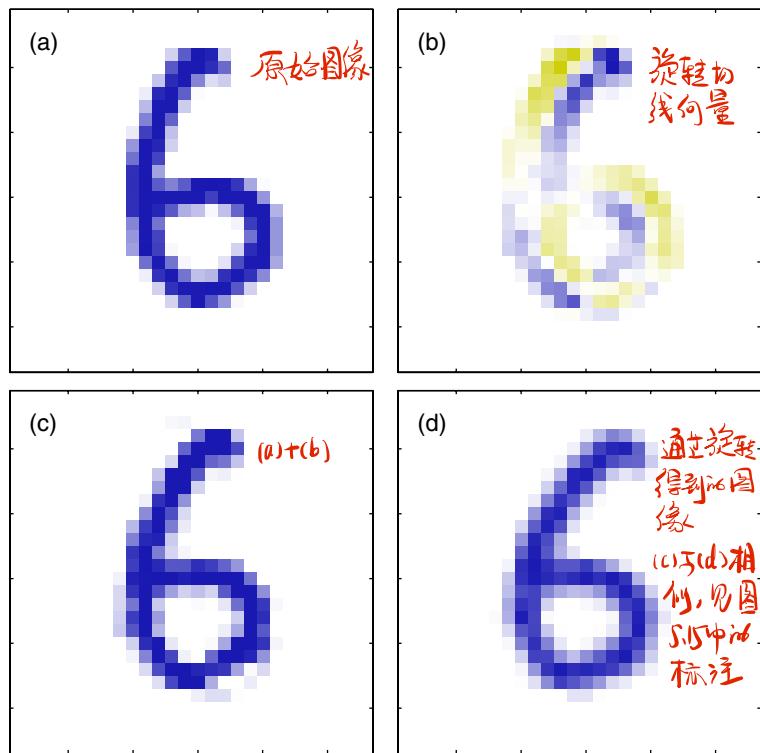
The regularization function depends on the network weights through the Jacobian  $\mathbf{J}$ . A backpropagation formalism for computing the derivatives of the regularizer with respect to the network weights is easily obtained by extension of the techniques introduced in Section 5.3.

If the transformation is governed by  $L$  parameters (e.g.,  $L = 3$  for the case of translations combined with in-plane rotations in a two-dimensional image), then the manifold  $\mathcal{M}$  will have dimensionality  $L$ , and the corresponding regularizer is given by the sum of terms of the form (5.128), one for each transformation. If several transformations are considered at the same time, and the network mapping is made invariant to each separately, then it will be (locally) invariant to combinations of the transformations (Simard *et al.*, 1992).

### Exercise 5.26

where blue and yellow correspond to positive and negative values, respectively

**Figure 5.16** Illustration showing (a) the original image  $\mathbf{x}$  of a handwritten digit, (b) the tangent vector  $\tau$  corresponding to an infinitesimal clockwise rotation, (c) the result of adding a small contribution from the tangent vector to the original image giving  $\mathbf{x} + \epsilon\tau$  with  $\epsilon = 15$  degrees, and (d) the true image rotated for comparison.



A related technique, called *tangent distance*, can be used to build invariance properties into distance-based methods such as nearest-neighbour classifiers (Simard *et al.*, 1993).

### 5.5.5 Training with transformed data

We have seen that one way to encourage invariance of a model to a set of transformations is to expand the training set using transformed versions of the original input patterns. Here we show that this approach is closely related to the technique of tangent propagation (Bishop, 1995b; Leen, 1995).

[ As in Section 5.5.4, we shall consider a transformation governed by a single parameter  $\xi$  and described by the function  $s(\mathbf{x}, \xi)$ , with  $s(\mathbf{x}, 0) = \mathbf{x}$ . We shall also consider a sum-of-squares error function. The error function for untransformed inputs can be written (in the infinite data set limit) in the form

$$E = \frac{1}{2} \iint \{y(\mathbf{x}) - t\}^2 p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt \quad (5.129)$$

as discussed in Section 1.5.5. Here we have considered a network having a single output, in order to keep the notation uncluttered. If we now consider an infinite number of copies of each data point, each of which is perturbed by the transformation

in which the parameter  $\xi$  is drawn from a distribution  $p(\xi)$ , then the error function defined over this expanded data set can be written as

$$\tilde{E} = \frac{1}{2} \iint \{y(\mathbf{s}(\mathbf{x}, \xi)) - t\}^2 p(t|\mathbf{x}) p(\mathbf{x}) p(\xi) d\mathbf{x} dt d\xi. \quad (5.130)$$

We now assume that the distribution  $p(\xi)$  has zero mean with small variance, so that we are only considering small transformations of the original input vectors. We can then expand the transformation function as a Taylor series in powers of  $\xi$  to give

$$\begin{aligned} \mathbf{s}(\mathbf{x}, \xi) &= \mathbf{s}(\mathbf{x}, 0) + \xi \frac{\partial}{\partial \xi} \mathbf{s}(\mathbf{x}, \xi) \Big|_{\xi=0} + \frac{\xi^2}{2} \frac{\partial^2}{\partial \xi^2} \mathbf{s}(\mathbf{x}, \xi) \Big|_{\xi=0} + O(\xi^3) \\ &= \mathbf{x} + \xi \boldsymbol{\tau} + \frac{1}{2} \xi^2 \boldsymbol{\tau}' + O(\xi^3) \end{aligned}$$

where  $\boldsymbol{\tau}'$  denotes the second derivative of  $\mathbf{s}(\mathbf{x}, \xi)$  with respect to  $\xi$  evaluated at  $\xi = 0$ . This allows us to expand the model function to give

$$y(\mathbf{s}(\mathbf{x}, \xi)) = y(\mathbf{x}) + \xi \boldsymbol{\tau}^T \nabla y(\mathbf{x}) + \frac{\xi^2}{2} \left[ (\boldsymbol{\tau}')^T \nabla y(\mathbf{x}) + \boldsymbol{\tau}^T \nabla \nabla y(\mathbf{x}) \boldsymbol{\tau} \right] + O(\xi^3).$$

Substituting into the mean error function (5.130) and expanding, we then have

$$\begin{aligned} \tilde{E} &= \frac{1}{2} \iint \{y(\mathbf{x}) - t\}^2 p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt \\ &+ \mathbb{E}[\xi] \iint \{y(\mathbf{x}) - t\} \boldsymbol{\tau}^T \nabla y(\mathbf{x}) p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt \\ &+ \mathbb{E}[\xi^2] \frac{1}{2} \iint \left[ \{y(\mathbf{x}) - t\} \frac{1}{2} \left\{ (\boldsymbol{\tau}')^T \nabla y(\mathbf{x}) + \boldsymbol{\tau}^T \nabla \nabla y(\mathbf{x}) \boldsymbol{\tau} \right\} \right. \\ &\quad \left. + (\boldsymbol{\tau}^T \nabla y(\mathbf{x}))^2 \right] p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt + O(\xi^3). \end{aligned}$$

Because the distribution of transformations has zero mean we have  $\mathbb{E}[\xi] = 0$ . Also, we shall denote  $\mathbb{E}[\xi^2]$  by  $\lambda$ . Omitting terms of  $O(\xi^3)$ , the average error function then becomes

$$\tilde{E} = E + \lambda \Omega \quad (5.131)$$

where  $E$  is the original sum-of-squares error, and the regularization term  $\Omega$  takes the form

$$\begin{aligned} \Omega &= \frac{1}{2} \int \left[ \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\} \frac{1}{2} \left\{ (\boldsymbol{\tau}')^T \nabla y(\mathbf{x}) + \boldsymbol{\tau}^T \nabla \nabla y(\mathbf{x}) \boldsymbol{\tau} \right\} \right. \\ &\quad \left. + (\boldsymbol{\tau}^T \nabla y(\mathbf{x}))^2 \right] p(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (5.132)$$

in which we have performed the integration over  $t$ .

We can further simplify this regularization term as follows. In Section 1.5.5 we saw that the function that minimizes the sum-of-squares error is given by the conditional average  $\mathbb{E}[t|\mathbf{x}]$  of the target values  $t$ . From (5.131) we see that the regularized error will equal the unregularized sum-of-squares plus terms which are  $O(\xi^2)$ , and so the network function that minimizes the total error will have the form

$$y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] + O(\xi^2). \quad (5.133)$$

Thus, to leading order in  $\xi^2$ , the first term in the regularizer vanishes and we are left with

$$\Omega = \frac{1}{2} \int (\tau^T \nabla y(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} \quad (5.134)$$

which is equivalent to the tangent propagation regularizer (5.128). ] If we consider the special case in which the transformation of the inputs simply consists of the addition of random noise, so that  $\mathbf{x} \rightarrow \mathbf{x} + \xi$ , then the regularizer takes the form

$$\Omega = \frac{1}{2} \int \|\nabla y(\mathbf{x})\|^2 p(\mathbf{x}) d\mathbf{x} \quad (5.135)$$

*Exercise 5.27*

其效果类似于

对 input  $\mathbf{x}$  添加随机噪声  
而数据增强

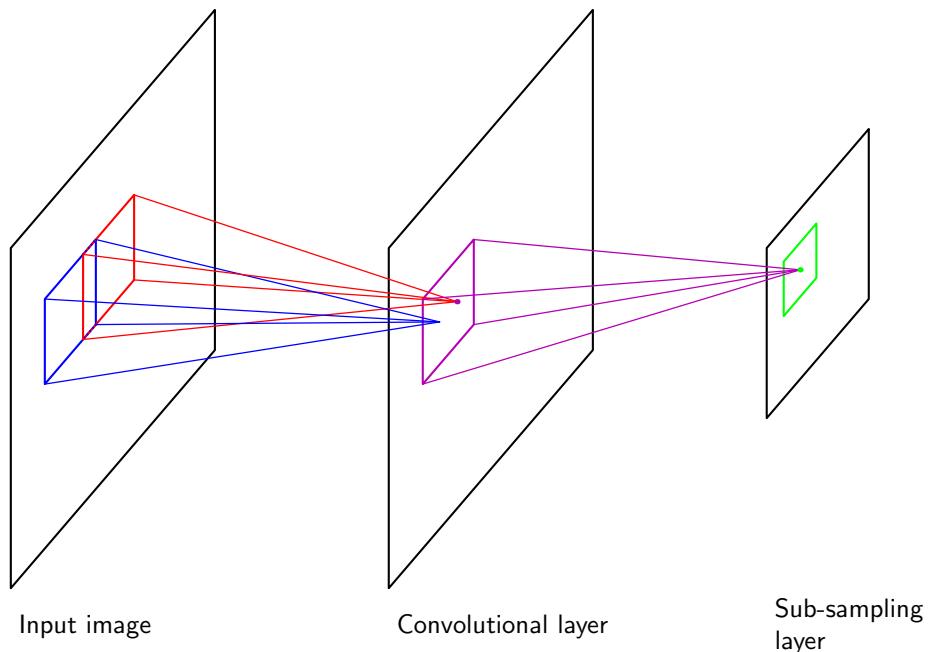
which is known as **Tikhonov regularization** (Tikhonov and Arsenin, 1977; Bishop, 1995b). Derivatives of this regularizer with respect to the network weights can be found using an extended backpropagation algorithm (Bishop, 1993). We see that, for small noise amplitudes, Tikhonov regularization is related to the addition of random noise to the inputs, which has been shown to improve generalization in appropriate circumstances (Sietsma and Dow, 1991). ]

## 5.5.6 Convolutional networks

Another approach to creating models that are invariant to certain transformations of the inputs is to build the invariance properties into the structure of a neural network. This is the basis for the **convolutional neural network** (Le Cun *et al.*, 1989; LeCun *et al.*, 1998), which has been widely applied to image data.

Consider the specific task of recognizing handwritten digits. Each input image comprises a set of pixel intensity values, and the desired output is a posterior probability distribution over the ten digit classes. We know that the identity of the digit is invariant under **translations** and **scaling** as well as (small) **rotations**. Furthermore, the network must also exhibit invariance to more subtle transformations such as **elastic deformations** of the kind illustrated in Figure 5.14. One simple approach would be to treat the image as the input to a fully connected network, such as the kind shown in Figure 5.1. Given a sufficiently large training set, such a network could in principle yield a good solution to this problem and would learn the appropriate invariances by example.

However, this approach ignores a key property of images, which is that nearby pixels are more strongly correlated than more distant pixels. Many of the modern approaches to computer vision exploit this property by extracting *local* features that depend only on small subregions of the image. Information from such features can then be merged in later stages of processing in order to detect higher-order features



**Figure 5.17** Diagram illustrating part of a convolutional neural network, showing a layer of convolutional units followed by a layer of subsampling units. Several successive pairs of such layers may be used.

and ultimately to yield information about the image as a whole. Also, local features that are useful in one region of the image are likely to be useful in other regions of the image, for instance if the object of interest is translated.

These notions are incorporated into convolutional neural networks through three mechanisms: (i) local receptive fields, (ii) weight sharing, and (iii) subsampling. The structure of a convolutional network is illustrated in Figure 5.17. In the convolutional layer the units are organized into planes, each of which is called a *feature map*. Units in a feature map each take inputs only from a small subregion of the image, and all of the units in a feature map are constrained to share the same weight values. For instance, a feature map might consist of 100 units arranged in a  $10 \times 10$  grid, with each unit taking inputs from a  $5 \times 5$  pixel patch of the image. The whole feature map therefore has 25 adjustable weight parameters plus one adjustable bias parameter. Input values from a patch are linearly combined using the weights and the bias, and the result transformed by a sigmoidal nonlinearity using (5.1). If we think of the units as feature detectors, then all of the units in a feature map detect the same pattern but at different locations in the input image. Due to the weight sharing, the evaluation of the activations of these units is equivalent to a convolution of the image pixel intensities with a ‘kernel’ comprising the weight parameters. If the input image is shifted, the activations of the feature map will be shifted by the same amount but will otherwise be unchanged. This provides the basis for the (approximate) invariance of

the network outputs to translations and distortions of the input image. Because we will typically need to detect multiple features in order to build an effective model, there will generally be multiple feature maps in the convolutional layer, each having its own set of weight and bias parameters.

The outputs of the convolutional units form the inputs to the subsampling layer of the network. For each feature map in the convolutional layer, there is a plane of units in the subsampling layer and each unit takes inputs from a small receptive field in the corresponding feature map of the convolutional layer. These units perform subsampling. For instance, each subsampling unit might take inputs from a  $2 \times 2$  unit region in the corresponding feature map and would compute the average of those inputs, multiplied by an adaptive weight with the addition of an adaptive bias parameter, and then transformed using a sigmoidal nonlinear activation function. The receptive fields are chosen to be contiguous and nonoverlapping so that there are half the number of rows and columns in the subsampling layer compared with the convolutional layer. In this way, the response of a unit in the subsampling layer will be relatively insensitive to small shifts of the image in the corresponding regions of the input space.

In a practical architecture, there may be several pairs of convolutional and subsampling layers. At each stage there is a larger degree of invariance to input transformations compared to the previous layer. There may be several feature maps in a given convolutional layer for each plane of units in the previous subsampling layer, so that the gradual reduction in spatial resolution is then compensated by an increasing number of features. The final layer of the network would typically be a fully connected, fully adaptive layer, with a softmax output nonlinearity in the case of multiclass classification.

The whole network can be trained by error minimization using backpropagation to evaluate the gradient of the error function. This involves a slight modification of the usual backpropagation algorithm to ensure that the shared-weight constraints are satisfied. Due to the use of local receptive fields, the number of weights in the network is smaller than if the network were fully connected. Furthermore, the number of independent parameters to be learned from the data is much smaller still, due to the substantial numbers of constraints on the weights.

### 5.5.7 Soft weight sharing

One way to reduce the effective complexity of a network with a large number of weights is to constrain weights within certain groups to be equal. This is the technique of weight sharing that was discussed in Section 5.5.6 as a way of building translation invariance into networks used for image interpretation. It is only applicable, however, to particular problems in which the form of the constraints can be specified in advance. Here we consider a form of **soft weight sharing** (Nowlan and Hinton, 1992) in which the hard constraint of equal weights is replaced by a form of regularization in which groups of weights are encouraged to have similar values. Furthermore, the division of weights into groups, the mean weight value for each group, and the spread of values within the groups are all determined as part of the learning process. }

#### Exercise 5.28

Page 270–272 修订:  
Section 5.5.7, from Equation (5.139) onwards: With the introduction of the  $\sigma_j^2$ , the regularization coefficient becomes irrelevant and hence it can be dropped from text and equations.

## 模型构建

### Section 2.3.9

Recall that the simple weight decay regularizer, given in (5.112), can be viewed as the negative log of a Gaussian prior distribution over the weights. We can encourage the weight values to form several groups, rather than just one group, by considering instead a probability distribution that is a *mixture* of Gaussians. The centres and variances of the Gaussian components, as well as the mixing coefficients, will be considered as adjustable parameters to be determined as part of the learning process. Thus, we have a probability density of the form

$$p(\mathbf{w}) = \prod_i p(w_i) \quad (5.136)$$

where

$$p(w_i) = \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \quad (5.137)$$

and  $\pi_j$  are the mixing coefficients. Taking the negative logarithm then leads to a regularization function of the form

$$\Omega(\mathbf{w}) = - \sum_i \ln \left( \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right). \quad (5.138)$$

The total error function is then given by

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \quad (5.139)$$

where  $\lambda$  is the regularization coefficient. This error is minimized both with respect to the weights  $w_i$  and with respect to the parameters  $\{\pi_j, \mu_j, \sigma_j\}$  of the mixture model. If the weights were constant, then the parameters of the mixture model could be determined by using the EM algorithm discussed in Chapter 9. However, the distribution of weights is itself evolving during the learning process, and so to avoid numerical instability, a joint optimization is performed simultaneously over the weights and the mixture-model parameters. This can be done using a standard optimization algorithm such as conjugate gradients or quasi-Newton methods.

In order to minimize the total error function, it is necessary to be able to evaluate its derivatives with respect to the various adjustable parameters. To do this it is convenient to regard the  $\{\pi_j\}$  as *prior* probabilities and to introduce the corresponding posterior probabilities which, following (2.192), are given by Bayes' theorem in the form

$$\gamma_j(w) = \frac{\pi_j \mathcal{N}(w | \mu_j, \sigma_j^2)}{\sum_k \pi_k \mathcal{N}(w | \mu_k, \sigma_k^2)}. \quad (5.140)$$

The derivatives of the total error function with respect to the weights are then given by

$$\frac{\partial \tilde{E}}{\partial w_i} = \frac{\partial E}{\partial w_i} + \lambda \sum_j \gamma_j(w_i) \frac{(w_i - \mu_j)}{\sigma_j^2}. \quad (5.141)$$

### Exercise 5.29

The effect of the regularization term is therefore to pull each weight towards the centre of the  $j^{\text{th}}$  Gaussian, with a force proportional to the posterior probability of that Gaussian for the given weight. This is precisely the kind of effect that we are seeking.

Derivatives of the error with respect to the centres of the Gaussians are also easily computed to give

$$\frac{\partial \tilde{E}}{\partial \mu_j} = \sum_i \gamma_j(w_i) \frac{(\mu_j - w_i)}{\sigma_j^2} \quad (5.142)$$

which has a simple intuitive interpretation, because it pushes  $\mu_j$  towards an average of the weight values, weighted by the posterior probabilities that the respective weight parameters were generated by component  $j$ . Similarly, the derivatives with respect to the variances are given by

$$\frac{\partial \tilde{E}}{\partial \sigma_j} = \sum_i \gamma_j(w_i) \left( \frac{1}{\sigma_j} - \frac{(w_i - \mu_j)^2}{\sigma_j^3} \right) \quad (5.143)$$

which drives  $\sigma_j$  towards the weighted average of the squared deviations of the weights around the corresponding centre  $\mu_j$ , where the weighting coefficients are again given by the posterior probability that each weight is generated by component  $j$ . Note that in a practical implementation, new variables  $\zeta_j$  defined by

$$\sigma_j^2 = \exp(\eta_j) \quad (5.144)$$

are introduced, and the minimization is performed with respect to the  $\zeta_j$ . This ensures that the parameters  $\sigma_j$  remain positive. It also has the effect of discouraging pathological solutions in which one or more of the  $\sigma_j$  goes to zero, corresponding to a Gaussian component collapsing onto one of the weight parameter values. Such solutions are discussed in more detail in the context of Gaussian mixture models in Section 9.2.1.

For the derivatives with respect to the mixing coefficients  $\pi_j$ , we need to take account of the constraints

$$\sum_j \pi_j = 1, \quad 0 \leq \pi_i \leq 1 \quad (5.145)$$

which follow from the interpretation of the  $\pi_j$  as prior probabilities. This can be done by expressing the mixing coefficients in terms of a set of auxiliary variables  $\{\eta_j\}$  using the *softmax* function given by

$$\pi_j = \frac{\exp(\eta_j)}{\sum_{k=1}^M \exp(\eta_k)}. \quad (5.146)$$

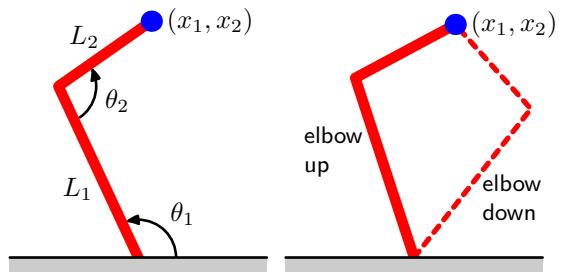
The derivatives of the regularized error function with respect to the  $\{\eta_j\}$  then take the form

*Exercise 5.30*

*Exercise 5.31*

*Exercise 5.32*

**Figure 5.18** The left figure shows a two-link robot arm, in which the Cartesian coordinates  $(x_1, x_2)$  of the end effector are determined uniquely by the two joint angles  $\theta_1$  and  $\theta_2$  and the (fixed) lengths  $L_1$  and  $L_2$  of the arms. This is known as the *forward kinematics* of the arm. In practice, we have to find the joint angles that will give rise to a desired end effector position and, as shown in the right figure, this *inverse kinematics* has two solutions corresponding to ‘elbow up’ and ‘elbow down’.



$$\frac{\partial \tilde{E}}{\partial \eta_j} = \sum_i \{\pi_j - \gamma_j(w_i)\}. \quad (5.147)$$

We see that  $\pi_j$  is therefore driven towards the average posterior probability for component  $j$ . } ]

## 5.6. Mixture Density Networks

The goal of supervised learning is to model a conditional distribution  $p(\mathbf{t}|\mathbf{x})$ , which for many simple regression problems is chosen to be Gaussian. However, practical machine learning problems can often have significantly non-Gaussian distributions. These can arise, for example, with *inverse problems* in which the distribution can be multimodal, in which case the Gaussian assumption can lead to very poor predictions.

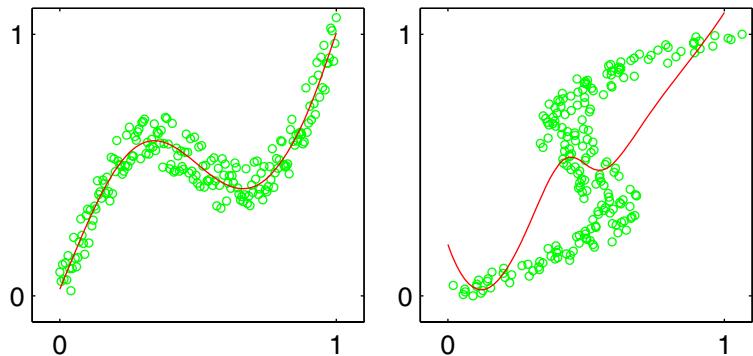
*Exercise 5.33*

As a simple example of an inverse problem, consider the kinematics of a robot arm, as illustrated in Figure 5.18. The *forward problem* involves finding the end effector position given the joint angles and has a unique solution. However, in practice we wish to move the end effector of the robot to a specific position, and to do this we must set appropriate joint angles. We therefore need to solve the inverse problem, which has two solutions as seen in Figure 5.18.

Forward problems often correspond to causality in a physical system and generally have a unique solution. For instance, a specific pattern of symptoms in the human body may be caused by the presence of a particular disease. In pattern recognition, however, we typically have to solve an inverse problem, such as trying to predict the presence of a disease given a set of symptoms. If the forward problem involves a many-to-one mapping, then the inverse problem will have multiple solutions. For instance, several different diseases may result in the same symptoms.

In the robotics example, the kinematics is defined by geometrical equations, and the multimodality is readily apparent. However, in many machine learning problems the presence of multimodality, particularly in problems involving spaces of high dimensionality, can be less obvious. For tutorial purposes, however, we shall consider a simple toy problem for which we can easily visualize the multimodality. Data for this problem is generated by sampling a variable  $x$  uniformly over the interval  $(0, 1)$ , to give a set of values  $\{x_n\}$ , and the corresponding target values  $t_n$  are obtained

**Figure 5.19** On the left is the data set for a simple ‘forward problem’ in which the red curve shows the result of fitting a two-layer neural network by minimizing the sum-of-squares error function. The corresponding inverse problem, shown on the right, is obtained by exchanging the roles of  $x$  and  $t$ . Here the same network trained again by minimizing the sum-of-squares error function gives a very poor fit to the data due to the multimodality of the data set.



by computing the function  $x_n + 0.3 \sin(2\pi x_n)$  and then adding uniform noise over the interval  $(-0.1, 0.1)$ . The inverse problem is then obtained by keeping the same data points but exchanging the roles of  $x$  and  $t$ . Figure 5.19 shows the data sets for the forward and inverse problems, along with the results of fitting two-layer neural networks having 6 hidden units and a single linear output unit by minimizing a sum-of-squares error function. Least squares corresponds to maximum likelihood under a Gaussian assumption. We see that this leads to a very poor model for the highly non-Gaussian inverse problem.

We therefore seek a general framework for modelling conditional probability distributions. This can be achieved by using a mixture model for  $p(t|x)$  in which both the mixing coefficients as well as the component densities are flexible functions of the input vector  $x$ , giving rise to the *mixture density network*. For any given value of  $x$ , the mixture model provides a general formalism for modelling an arbitrary conditional density function  $p(t|x)$ . Provided we consider a sufficiently flexible network, we then have a framework for approximating arbitrary conditional distributions.

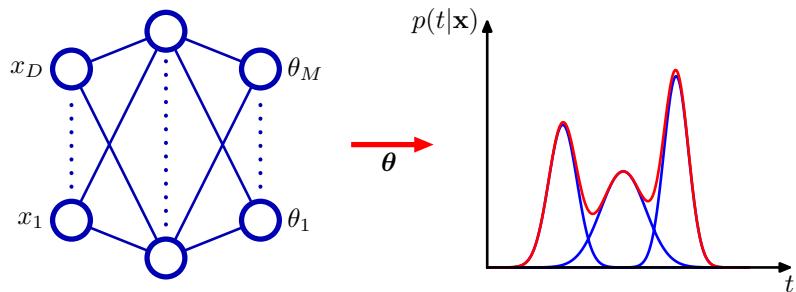
**举例** [ Here we shall develop the model explicitly for Gaussian components, so that

$$p(t|x) = \sum_{k=1}^K \pi_k(x) \mathcal{N}(t|\mu_k(x), \sigma_k^2(x)). \quad (5.148)$$

This is an example of a *heteroscedastic* model since the noise variance on the data is a function of the input vector  $x$ . Instead of Gaussians, we can use other distributions for the components, such as Bernoulli distributions if the target variables are binary rather than continuous. We have also specialized to the case of isotropic covariances for the components, although the mixture density network can readily be extended to allow for general covariance matrices by representing the covariances using a Cholesky factorization (Williams, 1996). Even with isotropic components, the conditional distribution  $p(t|x)$  does not assume factorization with respect to the components of  $t$  (in contrast to the standard sum-of-squares regression model) as a consequence of the mixture distribution.

这里说明是模型  $P(t|x)$  并未引介  
均独立性假设  
即  $P(t|x) = \prod P(t_i|x)$

We now take the various parameters of the mixture model, namely the mixing coefficients  $\pi_k(x)$ , the means  $\mu_k(x)$ , and the variances  $\sigma_k^2(x)$ , to be governed by



**Figure 5.20** The mixture density network can represent general conditional probability densities  $p(t|x)$  by considering a parametric mixture model for the distribution of  $t$  whose parameters are determined by the outputs of a neural network that takes  $x$  as its input vector.

the outputs of a conventional neural network that takes  $x$  as its input. The structure of this mixture density network is illustrated in Figure 5.20. The mixture density network is closely related to the mixture of experts discussed in Section 14.5.3. The principle difference is that in the mixture density network the same function is used to predict the parameters of all of the component densities as well as the mixing coefficients, and so the nonlinear hidden units are shared amongst the input-dependent functions.

**参数分析** The neural network in Figure 5.20 can, for example, be a two-layer network having sigmoidal ('tanh') hidden units. If there are  $K$  components in the mixture model (5.148), and if  $t$  has  $L$  components, then the network will have  $K$  output unit activations denoted by  $a_k^\pi$  that determine the mixing coefficients  $\pi_k(x)$ ,  $K$  outputs denoted by  $a_k^\sigma$  that determine the kernel widths  $\sigma_k(x)$ , and  $K \times L$  outputs denoted by  $a_{kj}^\mu$  that determine the components  $\mu_{kj}(x)$  of the kernel centres  $\mu_k(x)$ . The total number of network outputs is given by  $(K+2)L$ , as compared with the usual  $L$  outputs for a network, which simply predicts the conditional means of the target variables.

**模型参数化  
的具体情况** The mixing coefficients must satisfy the constraints

$$\sum_{k=1}^K \pi_k(x) = 1, \quad 0 \leq \pi_k(x) \leq 1 \quad (5.149)$$

which can be achieved using a set of softmax outputs

$$\pi_k(x) = \frac{\exp(a_k^\pi)}{\sum_{l=1}^K \exp(a_l^\pi)}. \quad (5.150)$$

Similarly, the variances must satisfy  $\sigma_k^2(x) \geq 0$  and so can be represented in terms of the exponentials of the corresponding network activations using

$$\sigma_k(x) = \exp(a_k^\sigma). \quad (5.151)$$

Finally, because the means  $\mu_k(x)$  have real components, they can be represented

directly by the network output activations

$$\mu_{kj}(\mathbf{x}) = a_{kj}^\mu. \quad (5.152)$$

**模型结构** The adaptive parameters of the mixture density network comprise the vector  $\mathbf{w}$  of weights and biases in the neural network, that can be set by maximum likelihood, or equivalently by minimizing an error function defined to be the negative logarithm of the likelihood. For independent data, this error function takes the form

$$E(\mathbf{w}) = - \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \right\} \quad (5.153)$$

where we have made the dependencies on  $\mathbf{w}$  explicit.

In order to minimize the error function, we need to calculate the derivatives of the error  $E(\mathbf{w})$  with respect to the components of  $\mathbf{w}$ . These can be evaluated by using the standard backpropagation procedure, provided we obtain suitable expressions for the derivatives of the error with respect to the output-unit activations. These represent error signals  $\delta$  for each pattern and for each output unit, and can be back-propagated to the hidden units and the error function derivatives evaluated in the usual way. Because the error function (5.153) is composed of a sum of terms, one for each training data point, we can consider the derivatives for a particular pattern  $n$  and then find the derivatives of  $E$  by summing over all patterns.

**模型推导** Because we are dealing with mixture distributions, it is convenient to view the mixing coefficients  $\pi_k(\mathbf{x})$  as  $\mathbf{x}$ -dependent prior probabilities and to introduce the corresponding posterior probabilities given by

$$\gamma_{nk} = \gamma_k(\mathbf{t}_n | \mathbf{x}_n) = \frac{\pi_k \mathcal{N}_{nk}}{\sum_{l=1}^K \pi_l \mathcal{N}_{nl}} \quad (5.154)$$

where  $\mathcal{N}_{nk}$  denotes  $\mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n), \sigma_k^2(\mathbf{x}_n))$ .

**Exercise 5.34** The derivatives with respect to the network output activations governing the mixing coefficients are given by

$$\frac{\partial E_n}{\partial a_k^\pi} = \pi_k - \gamma_{nk} \quad (5.155)$$

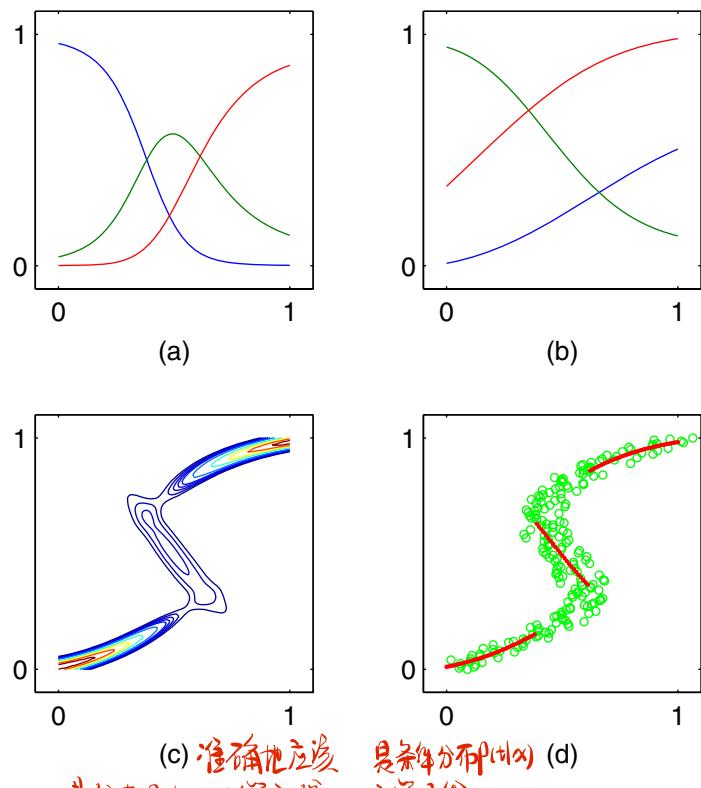
**Exercise 5.35** Similarly, the derivatives with respect to the output activations controlling the component means are given by

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = \gamma_{nk} \left\{ \frac{\mu_{kl} - \mathbf{t}_n}{\sigma_k^2} \right\}. \quad (5.156)$$

**Exercise 5.36** Finally, the derivatives with respect to the output activations controlling the component variances are given by

$$\frac{\partial E_n}{\partial a_k^\sigma} = -\gamma_{nk} \left\{ \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^2} - \frac{L}{\sigma_k^2} \right\}. \quad (5.157)$$

**Figure 5.21** (a) Plot of the mixing coefficients  $\pi_k(x)$  as a function of  $x$  for the three kernel functions in a mixture density network trained on the data shown in Figure 5.19. The model has three Gaussian components, and uses a two-layer multi-layer perceptron with five ‘tanh’ sigmoidal units in the hidden layer, and nine outputs (corresponding to the 3 means and 3 variances of the Gaussian components and the 3 mixing coefficients). At both small and large values of  $x$ , where the conditional probability density of the target data is unimodal, only one of the kernels has a high value for its prior probability, while at intermediate values of  $x$ , where the conditional density is trimodal, the three mixing coefficients have comparable values. (b) Plots of the means  $\mu_k(x)$  using the same colour coding as for the mixing coefficients. (c) Plot of the contours of the corresponding conditional probability density of the target data for the same mixture density network. (d) Plot of the approximate conditional mode, shown by the red points, of the conditional density.



年份的  
条件分布  
数  
是准确地应该  
是分布  $p(t|x)$  的等高线。  
b 不过  $P(t,x) = P(x)P(t|x)$   
 $P(x)$  为均匀分布，因此我们认为

We illustrate the use of a mixture density network by returning to the toy example of an inverse problem shown in Figure 5.19. Plots of the mixing coefficients  $\pi_k(x)$ , the means  $\mu_k(x)$ , and the conditional density contours corresponding to  $p(t|x)$ , are shown in Figure 5.21. The outputs of the neural network, and hence the parameters in the mixture model, are necessarily continuous single-valued functions of the input variables. However, we see from Figure 5.21(c) that the model is able to produce a conditional density that is unimodal for some values of  $x$  and trimodal for other values by modulating the amplitudes of the mixing components  $\pi_k(x)$ .

条件分布和条件逆  
Once a mixture density network has been trained, it can predict the conditional density function of the target data for any given value of the input vector. This conditional density represents a complete description of the generator of the data, so far as the problem of predicting the value of the output vector is concerned. From this density function we can calculate more specific quantities that may be of interest in different applications. One of the simplest of these is the mean, corresponding to the conditional average of the target data, and is given by

$$\mathbb{E}[\mathbf{t}|\mathbf{x}] = \int \mathbf{t} p(\mathbf{t}|\mathbf{x}) d\mathbf{t} = \sum_{k=1}^K \pi_k(\mathbf{x}) \mu_k(\mathbf{x}) \quad (5.158)$$

where we have used (5.148). Because a standard network trained by least squares is approximating the conditional mean, we see that a mixture density network can reproduce the conventional least-squares result as a special case. Of course, as we have already noted, for a multimodal distribution the conditional mean is of limited value.

We can similarly evaluate the variance of the density function about the conditional average, to give

$$s^2(\mathbf{x}) = \mathbb{E} [\|\mathbf{t} - \mathbb{E}[\mathbf{t}|\mathbf{x}]\|^2 | \mathbf{x}] \quad (5.159)$$

$$= \sum_{k=1}^K \pi_k(\mathbf{x}) \left\{ \sigma_k^2(\mathbf{x}) + \left\| \boldsymbol{\mu}_k(\mathbf{x}) - \sum_{l=1}^K \pi_l(\mathbf{x}) \boldsymbol{\mu}_l(\mathbf{x}) \right\|^2 \right\} \quad (5.160)$$

where we have used (5.148) and (5.158). This is more general than the corresponding least-squares result because the variance is a function of  $\mathbf{x}$ .]

**对模态分布，在预测时，使用条件均值作为结果输出，效果并不好，比如→此时条件众数 (conditional mode) 可能更有参考价值，但 mixture density network 的条件众数并不在简单地解析解，需要进行数值迭代计算，因此通常的替代方案是取混合模型中权重最大的 Component 的分布的均值作为多模态分布 (不是一个函数，one to many)**

## 5.7. Bayesian Neural Networks

So far, our discussion of neural networks has focussed on the use of maximum likelihood to determine the network parameters (weights and biases). Regularized maximum likelihood can be interpreted as a MAP (maximum posterior) approach in which the regularizer can be viewed as the logarithm of a prior parameter distribution. However, in a Bayesian treatment we need to marginalize over the distribution of parameters in order to make predictions.

In Section 3.3, we developed a Bayesian solution for a simple linear regression model under the assumption of Gaussian noise. We saw that the posterior distribution, which is Gaussian, could be evaluated exactly and that the predictive distribution could also be found in closed form. In the case of a multilayered network, the highly nonlinear dependence of the network function on the parameter values means that an exact Bayesian treatment can no longer be found. In fact, the log of the posterior distribution will be nonconvex, corresponding to the multiple local minima in the error function.

**方法1：**The technique of variational inference, to be discussed in Chapter 10, has been applied to Bayesian neural networks using a factorized Gaussian approximation

to the posterior distribution (Hinton and van Camp, 1993) and also using a full-covariance Gaussian (Barber and Bishop, 1998a; Barber and Bishop, 1998b). The most complete treatment, however, has been based on the Laplace approximation (MacKay, 1992c; MacKay, 1992b) and forms the basis for the discussion given here.

方法2，  
拉普拉斯  
积分法。

本书内容

We will approximate the posterior distribution by a Gaussian, centred at a mode of the true posterior. Furthermore, we shall assume that the covariance of this Gaussian is small so that the network function is approximately linear with respect to the parameters over the region of parameter space for which the posterior probability is significantly nonzero. With these two approximations, we will obtain models that are analogous to the linear regression and classification models discussed in earlier chapters and so we can exploit the results obtained there. We can then make use of the evidence framework to provide point estimates for the hyperparameters and to compare alternative models (for example, networks having different numbers of hidden units). To start with, we shall discuss the regression case and then later consider the modifications needed for solving classification tasks.

### 5.7.1 Posterior parameter distribution

模型描述

[ Consider the problem of predicting a single continuous target variable  $t$  from a vector  $\mathbf{x}$  of inputs (the extension to multiple targets is straightforward). We shall suppose that the conditional distribution  $p(t|\mathbf{x})$  is Gaussian, with an  $\mathbf{x}$ -dependent mean given by the output of a neural network model  $y(\mathbf{x}, \mathbf{w})$ , and with precision (inverse variance)  $\beta$

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}). \quad (5.161)$$

Similarly, we shall choose a prior distribution over the weights  $\mathbf{w}$  that is Gaussian of the form

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}). \quad (5.162)$$

For an i.i.d. data set of  $N$  observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , with a corresponding set of target values  $\mathcal{D} = \{t_1, \dots, t_N\}$ , the likelihood function is given by

$$p(\mathcal{D}|\mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|y(\mathbf{x}_n, \mathbf{w}), \beta^{-1}) \quad (5.163)$$

and so the resulting posterior distribution is then

$$p(\mathbf{w}|\mathcal{D}, \alpha, \beta) \propto p(\mathbf{w}|\alpha)p(\mathcal{D}|\mathbf{w}, \beta). \quad (5.164)$$

which, as a consequence of the nonlinear dependence of  $y(\mathbf{x}, \mathbf{w})$  on  $\mathbf{w}$ , will be non-Gaussian. ]

**问题处理** [ We can find a Gaussian approximation to the posterior distribution by using the Laplace approximation. To do this, we must first find a (local) maximum of the posterior, and this must be done using (iterative numerical optimization). As usual, it is convenient to maximize the logarithm of the posterior, which can be written in the

也就是梯度下降法求解

form

$$\ln p(\mathbf{w}|\mathcal{D}) = -\frac{\alpha}{2}\mathbf{w}^T\mathbf{w} - \frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \text{const} \quad (5.165)$$

which corresponds to a regularized sum-of-squares error function. Assuming for the moment that  $\alpha$  and  $\beta$  are fixed, we can find a maximum of the posterior, which we denote  $\mathbf{w}_{\text{MAP}}$ , by standard nonlinear optimization algorithms such as conjugate gradients, using error backpropagation to evaluate the required derivatives.

Having found a mode  $\mathbf{w}_{\text{MAP}}$ , we can then build a local Gaussian approximation by evaluating the matrix of second derivatives of the negative log posterior distribution. From (5.165), this is given by

$$\mathbf{A} = -\nabla\nabla \ln p(\mathbf{w}|\mathcal{D}, \alpha, \beta) = \alpha\mathbf{I} + \beta\mathbf{H} \quad (5.166)$$

where  $\mathbf{H}$  is the Hessian matrix comprising the second derivatives of the sum-of-squares error function with respect to the components of  $\mathbf{w}$ . Algorithms for computing and approximating the Hessian were discussed in Section 5.4. The corresponding Gaussian approximation to the posterior is then given from (4.134) by

参数  $\mathbf{w}$  后验分布的高斯近似  $q(\mathbf{w}|\mathcal{D}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{A}^{-1})$ . (5.167)

前面通过泰勒展开局部近似得到式(5.167)。下面 // Similarly, the predictive distribution is obtained by marginalizing with respect to this posterior distribution

$$p(t|\mathbf{x}, \mathcal{D}) = \int p(t|\mathbf{x}, \mathbf{w})q(\mathbf{w}|\mathcal{D}) d\mathbf{w}. \quad (5.168)$$

为了计算后验分布  $p(\mathbf{w}|\mathcal{D})$  式(5.168), 我们作近似式(5.169), 最终可得式(5.172).

However, even with the Gaussian approximation to the posterior, this integration is still analytically intractable due to the nonlinearity of the network function  $y(\mathbf{x}, \mathbf{w})$  as a function of  $\mathbf{w}$ . To make progress, we now assume that the posterior distribution has small variance compared with the characteristic scales of  $\mathbf{w}$  over which  $y(\mathbf{x}, \mathbf{w})$  is varying. This allows us to make a Taylor series expansion of the network function around  $\mathbf{w}_{\text{MAP}}$  and retain only the linear terms

$$y(\mathbf{x}, \mathbf{w}) \simeq y(\mathbf{x}, \mathbf{w}_{\text{MAP}}) + \mathbf{g}^T(\mathbf{w} - \mathbf{w}_{\text{MAP}}) \quad (5.169)$$

where we have defined

$$\mathbf{g} = \nabla_{\mathbf{w}} y(\mathbf{x}, \mathbf{w})|_{\mathbf{w}=\mathbf{w}_{\text{MAP}}}. \quad \text{见式(5.167)}$$

With this approximation, we now have a linear-Gaussian model with a Gaussian distribution for  $p(\mathbf{w})$  and a Gaussian for  $p(t|\mathbf{w})$  whose mean is a linear function of  $\mathbf{w}$  of the form

$$p(t|\mathbf{x}, \mathbf{w}, \beta) \simeq \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}_{\text{MAP}}) + \mathbf{g}^T(\mathbf{w} - \mathbf{w}_{\text{MAP}}), \beta^{-1}). \quad (5.171)$$

### Exercise 5.38

We can therefore make use of the general result (2.115) for the marginal  $p(t)$  to give

Bayesian predictive  $p(t|\mathbf{x}, \mathcal{D}, \alpha, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}_{\text{MAP}}), \sigma^2(\mathbf{x}))$  (5.172)

where the input-dependent variance is given by

$$\sigma^2(\mathbf{x}) = \beta^{-1} + \mathbf{g}^T \mathbf{A}^{-1} \mathbf{g}. \quad (5.173)$$

{ We see that the predictive distribution  $p(t|\mathbf{x}, \mathcal{D})$  is a Gaussian whose mean is given by the network function  $y(\mathbf{x}, \mathbf{w}_{\text{MAP}})$  with the parameter set to their MAP value. The variance has two terms, the first of which arises from the intrinsic noise on the target variable, whereas the second is an  $\mathbf{x}$ -dependent term that expresses the uncertainty in the interpolant due to the uncertainty in the model parameters  $\mathbf{w}$ . This should be compared with the corresponding predictive distribution for the linear regression model, given by (3.58) and (3.59). } ]

### 5.7.2 Hyperparameter optimization

超参数  $\alpha, \beta$   
marginal likelihood: max  
model evidence

[ So far, we have assumed that the hyperparameters  $\alpha$  and  $\beta$  are fixed and known. We can make use of the evidence framework, discussed in Section 3.5, together with the Gaussian approximation to the posterior obtained using the Laplace approximation, to obtain a practical procedure for choosing the values of such hyperparameters. ]

The marginal likelihood, or evidence, for the hyperparameters is obtained by integrating over the network weights

$$p(\mathcal{D}|\alpha, \beta) = \int p(\mathcal{D}|\mathbf{w}, \beta) p(\mathbf{w}|\alpha) d\mathbf{w}. \quad (5.174)$$

Exercise 5.39

This is easily evaluated by making use of the Laplace approximation result (4.135). Taking logarithms then gives

$$\ln p(\mathcal{D}|\alpha, \beta) \simeq -E(\mathbf{w}_{\text{MAP}}) - \frac{1}{2} \ln |\mathbf{A}| + \frac{W}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \quad (5.175)$$

where  $W$  is the total number of parameters in  $\mathbf{w}$ , and the regularized error function is defined by

$$E(\mathbf{w}_{\text{MAP}}) = \frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}_{\text{MAP}}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}_{\text{MAP}}^T \mathbf{w}_{\text{MAP}}. \quad (5.176)$$

We see that this takes the same form as the corresponding result (3.86) for the linear regression model.

In the evidence framework, we make point estimates for  $\alpha$  and  $\beta$  by maximizing  $\ln p(\mathcal{D}|\alpha, \beta)$ . Consider first the maximization with respect to  $\alpha$ , which can be done by analogy with the linear regression case discussed in Section 3.5.2. We first define the eigenvalue equation

$$\beta \mathbf{H} \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad (5.177)$$

where  $\mathbf{H}$  is the Hessian matrix comprising the second derivatives of the sum-of-squares error function, evaluated at  $\mathbf{w} = \mathbf{w}_{\text{MAP}}$ . By analogy with (3.92), we obtain

$$\alpha = \frac{\gamma}{\mathbf{w}_{\text{MAP}}^T \mathbf{w}_{\text{MAP}}} \quad (5.178)$$

**Section 3.5.3**

where  $\gamma$  represents the effective number of parameters and is defined by

$$\gamma = \sum_{i=1}^W \frac{\lambda_i}{\alpha + \lambda_i}. \quad (5.179)$$

Note that this result was exact for the linear regression case. For the nonlinear neural network, however, it ignores the fact that changes in  $\alpha$  will cause changes in the Hessian  $\mathbf{H}$ , which in turn will change the eigenvalues. We have therefore implicitly ignored terms involving the derivatives of  $\lambda_i$  with respect to  $\alpha$ .

Similarly, from (3.95) we see that maximizing the evidence with respect to  $\beta$  gives the re-estimation formula

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}_{\text{MAP}}) - t_n\}^2. \quad (5.180)$$

As with the linear model, we need to alternate between re-estimation of the hyperparameters  $\alpha$  and  $\beta$  and updating of the posterior distribution. The situation with a neural network model is more complex, however, due to the multimodality of the posterior distribution. As a consequence, the solution for  $\mathbf{w}_{\text{MAP}}$  found by maximizing the log posterior will depend on the initialization of  $\mathbf{w}$ . Solutions that differ only as a consequence of the interchange and sign reversal symmetries in the hidden units are identical so far as predictions are concerned, and it is irrelevant which of the equivalent solutions is found. However, there may be inequivalent solutions as well, and these will generally yield different values for the optimized hyperparameters.<sup>3]</sup>

[ In order to compare different models, for example neural networks having different numbers of hidden units, we need to evaluate the model evidence  $p(\mathcal{D})$ . This can be approximated by taking (5.175) and substituting the values of  $\alpha$  and  $\beta$  obtained from the iterative optimization of these hyperparameters. A more careful evaluation is obtained by marginalizing over  $\alpha$  and  $\beta$ , again by making a Gaussian approximation (MacKay, 1992c; Bishop, 1995a). In either case, it is necessary to evaluate the determinant  $|\mathbf{A}|$  of the Hessian matrix. This can be problematic in practice because the determinant, unlike the trace, is sensitive to the small eigenvalues that are often difficult to determine accurately. ]

The Laplace approximation is based on a local quadratic expansion around a mode of the posterior distribution over weights. We have seen in Section 5.1.1 that any given mode in a two-layer network is a member of a set of  $M!2^M$  equivalent modes that differ by interchange and sign-change symmetries, where  $M$  is the number of hidden units. When comparing networks having different numbers of hidden units, this can be taken into account by multiplying the evidence by a factor of  $M!2^M$ . ]

### 5.7.3 Bayesian neural networks for classification

So far, we have used the Laplace approximation to develop a Bayesian treatment of neural network regression models. We now discuss the modifications to

相当于存在  $M!2^M$   
个这样而使用 Laplace  
approximation 计算  
在局部区域，所以  
将一个局部区域计  
算的结果乘以  $M!2^M$ 。

这也说明对神经网络，Laplace approximation 是假  
设并不成立，Laplace approximation 实际上并不适用于此。

使用上面得到的  
Bayesian model evidence  
对不同网络结构的  
模型进行评估比较

this framework that arise when it is applied to classification. Here we shall consider a network having a single logistic sigmoid output corresponding to a two-class classification problem. The extension to networks with multiclass softmax outputs is straightforward. We shall build extensively on the analogous results for linear classification models discussed in Section 4.5, and so we encourage the reader to familiarize themselves with that material before studying this section.

**Exercise 5.40**

**模型描述** The log likelihood function for this model is given by

$$\ln p(\mathcal{D}|\mathbf{w}) = \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (5.181)$$

where  $t_n \in \{0, 1\}$  are the target values, and  $y_n \equiv y(\mathbf{x}_n, \mathbf{w})$ . Note that there is no hyperparameter  $\beta$ , because the data points are assumed to be correctly labelled. As before, the prior is taken to be an isotropic Gaussian of the form (5.162).

**问题处理** The first stage in applying the Laplace framework to this model is to initialize the hyperparameter  $\alpha$ , and then to determine the parameter vector  $\mathbf{w}$  by maximizing the log posterior distribution. This is equivalent to minimizing the regularized error function

$$E(\mathbf{w}) = -\ln p(\mathcal{D}|\mathbf{w}) + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \quad (5.182)$$

and can be achieved using error backpropagation combined with standard optimization algorithms, as discussed in Section 5.3.

Having found a solution  $\mathbf{w}_{MAP}$  for the weight vector, the next step is to evaluate the Hessian matrix  $\mathbf{H}$  comprising the second derivatives of the negative log likelihood function. This can be done, for instance, using the exact method of Section 5.4.5, or using the outer product approximation given by (5.85). The second derivatives of the negative log posterior can again be written in the form (5.166), and the Gaussian approximation to the posterior is then given by (5.167).

**Exercise 5.41**

// To optimize the hyperparameter  $\alpha$ , we again maximize the marginal likelihood, which is easily shown to take the form

$$\ln p(\mathcal{D}|\alpha) \simeq -E(\mathbf{w}_{MAP}) - \frac{1}{2} \ln |\mathbf{A}| + \frac{W}{2} \ln \alpha - \text{const} \quad (5.183)$$

where the regularized error function is defined by

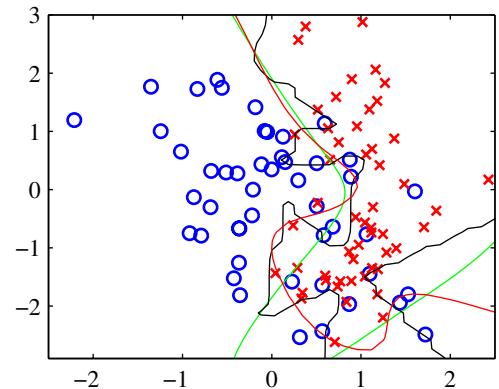
$$E(\mathbf{w}_{MAP}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + \frac{\alpha}{2} \mathbf{w}_{MAP}^T \mathbf{w}_{MAP} \quad (5.184)$$

in which  $y_n \equiv y(\mathbf{x}_n, \mathbf{w}_{MAP})$ . Maximizing this evidence function with respect to  $\alpha$  again leads to the re-estimation equation given by (5.178).

The use of the evidence procedure to determine  $\alpha$  is illustrated in Figure 5.22 for the synthetic two-dimensional data discussed in Appendix A.

// Finally, we need the predictive distribution, which is defined by (5.168). Again, this integration is intractable due to the nonlinearity of the network function. The

**Figure 5.22** Illustration of the evidence framework applied to a synthetic two-class data set. The green curve shows the optimal decision boundary, the black curve shows the result of fitting a two-layer network with 8 hidden units by maximum likelihood, and the red curve shows the result of including a regularizer in which  $\alpha$  is optimized using the evidence procedure, starting from the initial value  $\alpha = 0$ . Note that the evidence procedure greatly reduces the over-fitting of the network.



simplest approximation is to assume that the posterior distribution is very narrow and hence make the approximation

$$p(t|\mathbf{x}, \mathcal{D}) \simeq p(t|\mathbf{x}, \mathbf{w}_{\text{MAP}}). \quad (5.185)$$

**近似解二：**We can improve on this, however, by taking account of the variance of the posterior distribution. In this case, a linear approximation for the network outputs, as was used in the case of regression, would be inappropriate due to the logistic sigmoid output-unit activation function that constrains the output to lie in the range  $(0, 1)$ . Instead, we make a linear approximation for the output unit activation in the form

$$a(\mathbf{x}, \mathbf{w}) \simeq a_{\text{MAP}}(\mathbf{x}) + \mathbf{b}^T(\mathbf{w} - \mathbf{w}_{\text{MAP}}) \quad (5.186)$$

where  $a_{\text{MAP}}(\mathbf{x}) = a(\mathbf{x}, \mathbf{w}_{\text{MAP}})$ , and the vector  $\mathbf{b} \equiv \nabla a(\mathbf{x}, \mathbf{w}_{\text{MAP}})$  can be found by backpropagation.

Because we now have a Gaussian approximation for the posterior distribution over  $\mathbf{w}$ , and a model for  $a$  that is a linear function of  $\mathbf{w}$ , we can now appeal to the results of Section 4.5.2. The distribution of output unit activation values, induced by the distribution over network weights, is given by

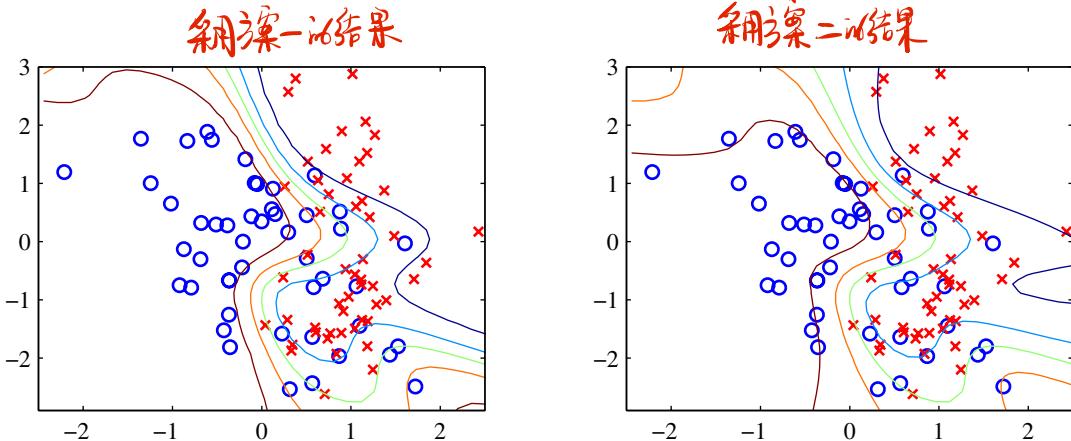
$$p(a|\mathbf{x}, \mathcal{D}) = \int \delta(a - a_{\text{MAP}}(\mathbf{x}) - \mathbf{b}^T(\mathbf{x})(\mathbf{w} - \mathbf{w}_{\text{MAP}})) q(\mathbf{w}|\mathcal{D}) d\mathbf{w} \quad (5.187)$$

where  $q(\mathbf{w}|\mathcal{D})$  is the Gaussian approximation to the posterior distribution given by (5.167). From Section 4.5.2, we see that this distribution is Gaussian with mean  $a_{\text{MAP}} \equiv a(\mathbf{x}, \mathbf{w}_{\text{MAP}})$ , and variance

$$\sigma_a^2(\mathbf{x}) = \mathbf{b}^T(\mathbf{x}) \mathbf{A}^{-1} \mathbf{b}(\mathbf{x}). \quad (5.188)$$

Finally, to obtain the predictive distribution, we must marginalize over  $a$  using

$$p(t=1|\mathbf{x}, \mathcal{D}) = \int \sigma(a) p(a|\mathbf{x}, \mathcal{D}) da. \quad (5.189)$$



**Figure 5.23** An illustration of the Laplace approximation for a Bayesian neural network having 8 hidden units with ‘tanh’ activation functions and a single logistic-sigmoid output unit. The weight parameters were found using scaled conjugate gradients, and the hyperparameter  $\alpha$  was optimized using the evidence framework. On the left is the result of using the simple approximation (5.185) based on a point estimate  $w_{MAP}$  of the parameters, in which the green curve shows the  $y = 0.5$  decision boundary, and the other contours correspond to output probabilities of  $y = 0.1, 0.3, 0.7$ , and  $0.9$ . On the right is the corresponding result obtained using (5.190). Note that the effect of marginalization is to spread out the contours and to make the predictions less confident, so that at each input point  $x$ , the posterior probabilities are shifted towards 0.5, while the  $y = 0.5$  contour itself is unaffected.

The convolution of a Gaussian with a logistic sigmoid is intractable. We therefore apply the approximation (4.153) to (5.189) giving

$$p(t=1|x, \mathcal{D}) = \sigma(\kappa(\sigma_a^2) b^T w_{MAP}) \quad (5.190)$$

where  $\kappa(\cdot)$  is defined by (4.154). Recall that both  $\sigma_a^2$  and  $b$  are functions of  $x$ . [3]

Figure 5.23 shows an example of this framework applied to the synthetic classification data set described in Appendix A.

## Exercises

- 5.1** (\*\*) Consider a two-layer network function of the form (5.7) in which the hidden-unit nonlinear activation functions  $h(\cdot)$  are given by logistic sigmoid functions of the form

$$\sigma(a) = \{1 + \exp(-a)\}^{-1}. \quad (5.191)$$

Show that there exists an equivalent network, which computes exactly the same function, but with hidden unit activation functions given by  $\tanh(a)$  where the  $\tanh$  function is defined by (5.59). Hint: first find the relation between  $\sigma(a)$  and  $\tanh(a)$ , and then show that the parameters of the two networks differ by linear transformations.

- 5.2** (\*) [www](#) Show that maximizing the likelihood function under the conditional distribution (5.16) for a multioutput neural network is equivalent to minimizing the sum-of-squares error function (5.11).

- 5.3** (\*\*) Consider a regression problem involving multiple target variables in which it is assumed that the distribution of the targets, conditioned on the input vector  $\mathbf{x}$ , is a Gaussian of the form

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{t}|\mathbf{y}(\mathbf{x}, \mathbf{w}), \Sigma) \quad (5.192)$$

where  $\mathbf{y}(\mathbf{x}, \mathbf{w})$  is the output of a neural network with input vector  $\mathbf{x}$  and weight vector  $\mathbf{w}$ , and  $\Sigma$  is the covariance of the assumed Gaussian noise on the targets. Given a set of independent observations of  $\mathbf{x}$  and  $\mathbf{t}$ , write down the error function that must be minimized in order to find the maximum likelihood solution for  $\mathbf{w}$ , if we assume that  $\Sigma$  is fixed and known. Now assume that  $\Sigma$  is also to be determined from the data, and write down an expression for the maximum likelihood solution for  $\Sigma$ . Note that the optimizations of  $\mathbf{w}$  and  $\Sigma$  are now coupled, in contrast to the case of independent target variables discussed in Section 5.2.

- 5.4** (\*\*) Consider a binary classification problem in which the target values are  $t \in \{0, 1\}$ , with a network output  $y(\mathbf{x}, \mathbf{w})$  that represents  $p(t = 1|\mathbf{x})$ , and suppose that there is a probability  $\epsilon$  that the class label on a training data point has been incorrectly set. Assuming independent and identically distributed data, write down the error function corresponding to the negative log likelihood. Verify that the error function (5.21) is obtained when  $\epsilon = 0$ . Note that this error function makes the model robust to incorrectly labelled data, in contrast to the usual error function.

- 5.5** (\*) **www** Show that maximizing likelihood for a multiclass neural network model in which the network outputs have the interpretation  $y_k(\mathbf{x}, \mathbf{w}) = p(t_k = 1|\mathbf{x})$  is equivalent to the minimization of the cross-entropy error function (5.24).

- 5.6** (\*) **www** Show the derivative of the error function (5.21) with respect to the activation  $a_k$  for an output unit having a logistic sigmoid activation function satisfies (5.18).

- 5.7** (\*) Show the derivative of the error function (5.24) with respect to the activation  $a_k$  for output units having a softmax activation function satisfies (5.18).

- 5.8** (\*) We saw in (4.88) that the derivative of the logistic sigmoid activation function can be expressed in terms of the function value itself. Derive the corresponding result for the ‘tanh’ activation function defined by (5.59).

- 5.9** (\*) **www** The error function (5.21) for binary classification problems was derived for a network having a logistic-sigmoid output activation function, so that  $0 \leq y(\mathbf{x}, \mathbf{w}) \leq 1$ , and data having target values  $t \in \{0, 1\}$ . Derive the corresponding error function if we consider a network having an output  $-1 \leq y(\mathbf{x}, \mathbf{w}) \leq 1$  and target values  $t = 1$  for class  $C_1$  and  $t = -1$  for class  $C_2$ . What would be the appropriate choice of output unit activation function?

- 5.10** (\*) **www** Consider a Hessian matrix  $\mathbf{H}$  with eigenvector equation (5.33). By setting the vector  $\mathbf{v}$  in (5.39) equal to each of the eigenvectors  $\mathbf{u}_i$  in turn, show that  $\mathbf{H}$  is positive definite if, and only if, all of its eigenvalues are positive.

- 5.11** (\*\*) **www** Consider a quadratic error function defined by (5.32), in which the Hessian matrix  $\mathbf{H}$  has an eigenvalue equation given by (5.33). Show that the contours of constant error are ellipses whose axes are aligned with the eigenvectors  $\mathbf{u}_i$ , with lengths that are inversely proportional to the square root of the corresponding eigenvalues  $\lambda_i$ .
- 5.12** (\*\*) **www** By considering the local Taylor expansion (5.32) of an error function about a stationary point  $\mathbf{w}^*$ , show that the necessary and sufficient condition for the stationary point to be a local minimum of the error function is that the Hessian matrix  $\mathbf{H}$ , defined by (5.30) with  $\widehat{\mathbf{w}} = \mathbf{w}^*$ , be positive definite.
- 5.13** (\*) Show that as a consequence of the symmetry of the Hessian matrix  $\mathbf{H}$ , the number of independent elements in the quadratic error function (5.28) is given by  $W(W + 3)/2$ .
- 5.14** (\*) By making a Taylor expansion, verify that the terms that are  $O(\epsilon)$  cancel on the right-hand side of (5.69).
- 5.15** (\*\*) In Section 5.3.4, we derived a procedure for evaluating the Jacobian matrix of a neural network using a backpropagation procedure. Derive an alternative formalism for finding the Jacobian based on *forward propagation* equations.
- 5.16** (\*) The outer product approximation to the Hessian matrix for a neural network using a sum-of-squares error function is given by (5.84). Extend this result to the case of multiple outputs.
- 5.17** (\*) Consider a squared loss function of the form

$$E = \frac{1}{2} \iint \{y(\mathbf{x}, \mathbf{w}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt \quad (5.193)$$

where  $y(\mathbf{x}, \mathbf{w})$  is a parametric function such as a neural network. The result (1.89) shows that the function  $y(\mathbf{x}, \mathbf{w})$  that minimizes this error is given by the conditional expectation of  $t$  given  $\mathbf{x}$ . Use this result to show that the second derivative of  $E$  with respect to two elements  $w_r$  and  $w_s$  of the vector  $\mathbf{w}$ , is given by

$$\frac{\partial^2 E}{\partial w_r \partial w_s} = \int \frac{\partial y}{\partial w_r} \frac{\partial y}{\partial w_s} p(\mathbf{x}) \, d\mathbf{x}. \quad (5.194)$$

Note that, for a finite sample from  $p(\mathbf{x})$ , we obtain (5.84).

- 5.18** (\*) Consider a two-layer network of the form shown in Figure 5.1 with the addition of extra parameters corresponding to skip-layer connections that go directly from the inputs to the outputs. By extending the discussion of Section 5.3.2, write down the equations for the derivatives of the error function with respect to these additional parameters.
- 5.19** (\*) **www** Derive the expression (5.85) for the outer product approximation to the Hessian matrix for a network having a single output with a logistic sigmoid output-unit activation function and a cross-entropy error function, corresponding to the result (5.84) for the sum-of-squares error function.

修订 Exercise 5.21: The text in the exercise could be misunderstood; a less ambiguous formulation is: “Extend the expression (5.86) for the outer product approximation of the Hessian matrix to the case of  $K > 1$  output units. Hence, derive a form that allows (5.87) to be used to incorporate sequentially contributions from individual outputs as well as individual patterns. This, together with the identity (5.88), will allow the use of (5.89) for finding the inverse of the Hessian by sequentially incorporating contributions from individual outputs and patterns.”

- 5.20** (★) Derive an expression for the outer product approximation to the Hessian matrix for a network having  $K$  outputs with a softmax output-unit activation function and a cross-entropy error function, corresponding to the result (5.84) for the sum-of-squares error function.
- 5.21** (★★★) Extend the expression (5.86) for the outer product approximation of the Hessian matrix to the case of  $K > 1$  output units. Hence, derive a recursive expression analogous to (5.87) for incrementing the number  $N$  of patterns and a similar expression for incrementing the number  $K$  of outputs. Use these results, together with the identity (5.88), to find sequential update expressions analogous to (5.89) for finding the inverse of the Hessian by incrementally including both extra patterns and extra outputs.
- 5.22** (★★) Derive the results (5.93), (5.94), and (5.95) for the elements of the Hessian matrix of a two-layer feed-forward network by application of the chain rule of calculus.
- 5.23** (★★) Extend the results of Section 5.4.5 for the exact Hessian of a two-layer network to include skip-layer connections that go directly from inputs to outputs.
- 5.24** (★) Verify that the network function defined by (5.113) and (5.114) is invariant under the transformation (5.115) applied to the inputs, provided the weights and biases are simultaneously transformed using (5.116) and (5.117). Similarly, show that the network outputs can be transformed according (5.118) by applying the transformation (5.119) and (5.120) to the second-layer weights and biases.
- 5.25** (★★★) **www** Consider a quadratic error function of the form

$$E = E_0 + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \quad (5.195)$$

where  $\mathbf{w}^*$  represents the minimum, and the Hessian matrix  $\mathbf{H}$  is positive definite and constant. Suppose the initial weight vector  $\mathbf{w}^{(0)}$  is chosen to be at the origin and is updated using simple gradient descent

$$\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \rho \nabla E \quad (5.196)$$

where  $\tau$  denotes the step number, and  $\rho$  is the learning rate (which is assumed to be small). Show that, after  $\tau$  steps, the components of the weight vector parallel to the eigenvectors of  $\mathbf{H}$  can be written

$$w_j^{(\tau)} = \{1 - (1 - \rho \eta_j)^\tau\} w_j^* \quad (5.197)$$

where  $w_j = \mathbf{w}^T \mathbf{u}_j$ , and  $\mathbf{u}_j$  and  $\eta_j$  are the eigenvectors and eigenvalues, respectively, of  $\mathbf{H}$  so that

$$\mathbf{H}\mathbf{u}_j = \eta_j \mathbf{u}_j. \quad (5.198)$$

Show that as  $\tau \rightarrow \infty$ , this gives  $\mathbf{w}^{(\tau)} \rightarrow \mathbf{w}^*$  as expected, provided  $|1 - \rho \eta_j| < 1$ . Now suppose that training is halted after a finite number  $\tau$  of steps. Show that the

components of the weight vector parallel to the eigenvectors of the Hessian satisfy

$$w_j^{(\tau)} \simeq w_j^* \quad \text{when} \quad \eta_j \gg (\rho\tau)^{-1} \quad (5.199)$$

$$|w_j^{(\tau)}| \ll |w_j^*| \quad \text{when} \quad \eta_j \ll (\rho\tau)^{-1}. \quad (5.200)$$

Compare this result with the discussion in Section 3.5.3 of regularization with simple weight decay, and hence show that  $(\rho\tau)^{-1}$  is analogous to the regularization parameter  $\lambda$ . The above results also show that the effective number of parameters in the network, as defined by (3.91), grows as the training progresses.

- 5.26** (\*\*) Consider a multilayer perceptron with arbitrary feed-forward topology, which is to be trained by minimizing the *tangent propagation* error function (5.127) in which the regularizing function is given by (5.128). Show that the regularization term  $\Omega$  can be written as a sum over patterns of terms of the form

$$\Omega_n = \frac{1}{2} \sum_k (\mathcal{G}y_k)^2 \quad (5.201)$$

where  $\mathcal{G}$  is a differential operator defined by

$$\mathcal{G} \equiv \sum_i \tau_i \frac{\partial}{\partial x_i}. \quad (5.202)$$

By acting on the forward propagation equations

$$z_j = h(a_j), \quad a_j = \sum_i w_{ji} z_i \quad (5.203)$$

with the operator  $\mathcal{G}$ , show that  $\Omega_n$  can be evaluated by forward propagation using the following equations:

$$\alpha_j = h'(a_j)\beta_j, \quad \beta_j = \sum_i w_{ji}\alpha_i. \quad (5.204)$$

where we have defined the new variables

$$\alpha_j \equiv \mathcal{G}z_j, \quad \beta_j \equiv \mathcal{G}a_j. \quad (5.205)$$

Now show that the derivatives of  $\Omega_n$  with respect to a weight  $w_{rs}$  in the network can be written in the form

$$\frac{\partial \Omega_n}{\partial w_{rs}} = \sum_k \alpha_k \{ \phi_{kr} z_s + \delta_{kr} \alpha_s \} \quad (5.206)$$

where we have defined

$$\delta_{kr} \equiv \frac{\partial y_k}{\partial a_r}, \quad \phi_{kr} \equiv \mathcal{G}\delta_{kr}. \quad (5.207)$$

Write down the backpropagation equations for  $\delta_{kr}$ , and hence derive a set of back-propagation equations for the evaluation of the  $\phi_{kr}$ .

- 5.27** (\*\*) **www** Consider the framework for training with transformed data in the special case in which the transformation consists simply of the addition of random noise  $\mathbf{x} \rightarrow \mathbf{x} + \boldsymbol{\xi}$  where  $\boldsymbol{\xi}$  has a Gaussian distribution with zero mean and unit covariance. By following an argument analogous to that of Section 5.5.5, show that the resulting regularizer reduces to the Tikhonov form (5.135).
- 5.28** (\*) **www** Consider a neural network, such as the convolutional network discussed in Section 5.5.6, in which multiple weights are constrained to have the same value. Discuss how the standard backpropagation algorithm must be modified in order to ensure that such constraints are satisfied when evaluating the derivatives of an error function with respect to the adjustable parameters in the network.
- 5.29** (\*) **www** Verify the result (5.141).
- 5.30** (\*) Verify the result (5.142).
- 5.31** (\*) Verify the result (5.143).
- 5.32** (\*\*) Show that the derivatives of the mixing coefficients  $\{\pi_k\}$ , defined by (5.146), with respect to the auxiliary parameters  $\{\eta_j\}$  are given by

$$\frac{\partial \pi_k}{\partial \eta_j} = \delta_{jk} \pi_j - \pi_j \pi_k. \quad (5.208)$$

$\sum_k \pi_k = 1 \text{ for all } i$

Hence, by making use of the constraint  $\sum_k \pi_k = 1$ , derive the result (5.147).

- 5.33** (\*) Write down a pair of equations that express the Cartesian coordinates  $(x_1, x_2)$  for the robot arm shown in Figure 5.18 in terms of the joint angles  $\theta_1$  and  $\theta_2$  and the lengths  $L_1$  and  $L_2$  of the links. Assume the origin of the coordinate system is given by the attachment point of the lower arm. These equations define the ‘forward kinematics’ of the robot arm.
- 5.34** (\*) **www** Derive the result (5.155) for the derivative of the error function with respect to the network output activations controlling the mixing coefficients in the mixture density network.
- 5.35** (\*) Derive the result (5.156) for the derivative of the error function with respect to the network output activations controlling the component means in the mixture density network.
- 5.36** (\*) Derive the result (5.157) for the derivative of the error function with respect to the network output activations controlling the component variances in the mixture density network.
- 5.37** (\*) Verify the results (5.158) and (5.160) for the conditional mean and variance of the mixture density network model.
- 5.38** (\*) Using the general result (2.115), derive the predictive distribution (5.172) for the Laplace approximation to the Bayesian neural network model.

- 5.39** (\*) **www** Make use of the Laplace approximation result (4.135) to show that the evidence function for the hyperparameters  $\alpha$  and  $\beta$  in the Bayesian neural network model can be approximated by (5.175).
- 5.40** (\*) **www** Outline the modifications needed to the framework for Bayesian neural networks, discussed in Section 5.7.3, to handle multiclass problems using networks having softmax output-unit activation functions. *and 5.7.2*
- 5.41** (\*\*) By following analogous steps to those given in Section 5.7.1 for regression networks, derive the result (5.183) for the marginal likelihood in the case of a network having a cross-entropy error function and logistic-sigmoid output-unit activation function.