

## SINGLE-LAYER NETWORKS

In Chapter 1 we showed that the optimal decision rule for minimizing the probability of misclassification requires a new pattern to be assigned to the class having the largest posterior probability. We also showed how the posterior probabilities can be related to class-conditional densities through Bayes' theorem, and in Chapter 2 we described several techniques for estimating these densities. An alternative approach, which circumvents the determination of probability densities, is based on the idea of a discriminant function, also introduced in Chapter 1. In a practical application of discriminant functions, specific parametrized functional forms are chosen, and the values of the parameters are then determined from a set of training data by means of a suitable learning algorithm.

The simplest choice of discriminant function consists of a linear combination of the input variables, in which the coefficients in the linear combination are the parameters of the model, and has been considered widely in the literature on conventional approaches to pattern recognition. This simple discriminant can be generalized by transforming the linear combination with a non-linear function (called an activation function) which leads to concepts such as logistic regression and the perceptron. Another extension involves transforming the input variables with fixed non-linear functions before forming the linear combination, to give generalized linear discriminants. As we shall see, these various forms of linear discriminant can be regarded as forms of neural network in which there is a single layer of adaptive weights between the inputs and the outputs.

Various techniques exist for determining the weight values in single-layer networks, and in this chapter we shall consider several of them in detail. In particular, we shall study perceptron learning, least-squares methods and the Fisher discriminant. As well as forming an important class of techniques in their own right, single-layer networks provide many useful insights into the properties of more complex multi-layer networks. Single-layer networks were widely studied in the 1960's, and the history of such networks is reviewed in Widrow and Lehr (1990). Two useful books from this period are Nilsson (1965) and Lewis and Coates (1967).

### 3.1 Linear discriminant functions

In Chapter 1 we saw that optimal discriminant functions can be determined from class-conditional densities via Bayes' theorem. Instead of performing density estimation, however, we can postulate specific parametrized functional forms for

the discriminant functions and use the training data set to determine suitable values for the parameters. In this section we consider various forms of linear discriminant, and discuss their properties.

### 3.1.1 Two classes

We begin by considering the two-category classification problem. In Chapter 1 we introduced the concept of a discriminant function  $y(\mathbf{x})$  such that the vector  $\mathbf{x}$  is assigned to class  $C_1$  if  $y(\mathbf{x}) > 0$  and to class  $C_2$  if  $y(\mathbf{x}) < 0$ . The simplest choice of discriminant function is one which is linear in the components of  $\mathbf{x}$ , and which can therefore be written as

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (3.1)$$

where we shall refer to the  $d$ -dimensional vector  $\mathbf{w}$  as the *weight vector* and the parameter  $w_0$  as the *bias*. Sometimes  $-w_0$  is called a *threshold*. Note that the use of the term bias here is quite distinct from the concept of statistical bias which is discussed briefly on page 41, and at length in Section 9.1. From Section 2.1.3 we know that, for class-conditional densities having normal distributions with equal covariance matrices, a linear discriminant of the form (3.1) is optimal.

The expression in (3.1) has a simple geometrical interpretation (Duda and Hart, 1973) as follows. We first note that the decision boundary  $y(\mathbf{x}) = 0$  corresponds to a  $(d - 1)$ -dimensional hyperplane in  $d$ -dimensional  $\mathbf{x}$ -space. For the case of a two-dimensional input space,  $d = 2$ , the decision boundary is a straight line, as shown in Figure 3.1. If  $\mathbf{x}^A$  and  $\mathbf{x}^B$  are two points on the hyperplane, then  $y(\mathbf{x}^A) = 0 = y(\mathbf{x}^B)$  and so, using (3.1), we have  $\mathbf{w}^T(\mathbf{x}^B - \mathbf{x}^A) = 0$ . Thus,  $\mathbf{w}$  is normal to any vector lying in the hyperplane, and so we see that  $\mathbf{w}$  determines the orientation of the decision boundary. If  $\mathbf{x}$  is a point on the hyperplane then the normal distance from the origin to the hyperplane is given by

$$x_d = \frac{w_1 x_1 + \dots + w_{d-1} x_{d-1}}{w_d}$$

$$l = \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|} \quad (3.2)$$

where we have used  $y(\mathbf{x}) = 0$  together with (3.1). Thus, the bias  $w_0$  determines the position of the hyperplane in  $\mathbf{x}$ -space, as indicated in Figure 3.1.

There is a slightly different notation which we can adopt which will often prove convenient. If we define new  $(d + 1)$ -dimensional vectors  $\tilde{\mathbf{w}} = (w_0, \mathbf{w})$  and  $\tilde{\mathbf{x}} = (1, \mathbf{x})$ , then we can rewrite (3.1) in the form

$$y(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}. \quad (3.3)$$

With this notation we can interpret the decision boundary  $y(\mathbf{x}) = 0$  as a  $d$ -dimensional hyperplane which passes through the origin in  $(d + 1)$ -dimensional  $\tilde{\mathbf{x}}$ -space.

We can represent the linear discriminant function in (3.1) or (3.3) in terms

$$y = \alpha = w^T x + w_0$$

### 3.1: Linear discriminant functions

79

$$\Leftrightarrow w^T x = -w_0$$

dot product is scalar projection (wiki),  $\frac{w^T x}{\|w\|} = \text{projection on } x$

$$\Leftrightarrow \frac{-w_0}{\|w\|}$$

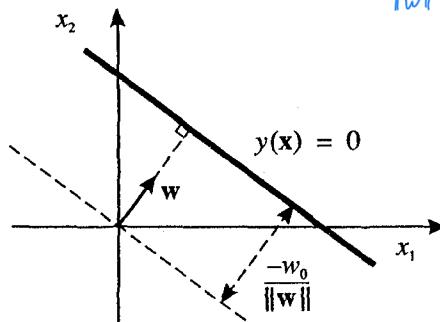


Figure 3.1. A linear decision boundary, corresponding to  $y(\mathbf{x}) = 0$ , in a two-dimensional input space  $(x_1, x_2)$ . The weight vector  $\mathbf{w}$ , which can be represented as a vector in  $\mathbf{x}$ -space, defines the orientation of the decision plane, while the bias  $w_0$  defines the position of the plane in terms of its perpendicular distance from the origin.

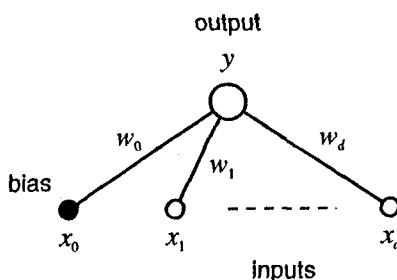


Figure 3.2. Representation of a linear discriminant function as a neural network diagram. Each component in the diagram corresponds to a variable in the linear discriminant expression. The bias  $w_0$  can be considered as a weight parameter from an extra input whose activation  $x_0$  is permanently set to +1.

of a network diagram as shown in Figure 3.2. Inputs  $x_1, \dots, x_d$  are shown as circles, which are connected by the weights  $w_1, \dots, w_d$  to the output  $y(\mathbf{x})$ . The bias  $w_0$  is represented as a weight from an extra input  $x_0$  which is permanently set to unity.

### 3.1.2 Several classes

Linear discriminants can easily be extended to the case of  $c$  classes by following the ideas introduced in Chapter 1 and using one discriminant function  $y_k(\mathbf{x})$  for each class  $\mathcal{C}_k$  of the form

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}. \quad (3.4)$$

A new point  $\mathbf{x}$  is then assigned to class  $\mathcal{C}_k$  if  $y_k(\mathbf{x}) > y_j(\mathbf{x})$  for all  $j \neq k$ . The decision boundary separating class  $\mathcal{C}_k$  from class  $\mathcal{C}_j$  is given by  $y_k(\mathbf{x}) = y_j(\mathbf{x})$  which, for linear discriminants, corresponds to a hyperplane of the form

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0. \quad (3.5)$$

By analogy with our earlier results for the single discriminant (3.1), we see that the normal to the decision boundary is given by the difference between the two weight vectors, and that the perpendicular distance of the decision boundary from the origin is given by

$$l = -\frac{(w_{k0} - w_{j0})}{\|\mathbf{w}_k - \mathbf{w}_j\|}. \quad (3.6)$$

The multiclass linear discriminant function (3.4) can be expressed in terms of a neural network diagram as shown in Figure 3.3. The circles at the top of the diagram, corresponding to the functions  $y_k(\mathbf{x})$  in (3.4) are sometimes called *processing units*, and the evaluation of the discriminant functions can be viewed as a flow of information from the inputs to the outputs. Each output  $y_k(\mathbf{x})$  is associated with a weight vector  $\mathbf{w}_k$  and a bias  $w_{k0}$ . We can express the network outputs in terms of the components of the vectors  $\{\mathbf{w}_k\}$  to give

$$y_k(\mathbf{x}) = \sum_{i=1}^d w_{ki} x_i + w_{k0}. \quad (3.7)$$

Then each line in Figure 3.3 connecting an input  $i$  to an output  $k$  corresponds to a weight parameter  $w_{ki}$ . As before, we can regard the bias parameters as being weights from an extra input  $x_0 = 1$ , so that

$$y_k(\mathbf{x}) = \sum_{i=0}^d w_{ki} x_i. \quad (3.8)$$

Once the network is trained, a new vector is classified by applying it to the inputs of the network, computing the output unit activations, and assigning the vector to the class whose output unit has the largest activation. This leads to a set of decision regions which are always simply connected and convex. To see

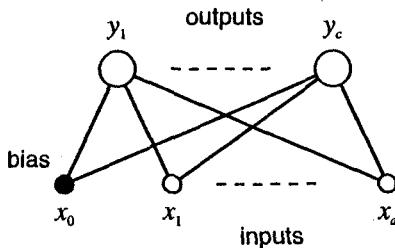


Figure 3.3. Representation of multiple linear discriminant functions  $y_k(\mathbf{x})$  as a neural network diagram having  $c$  output units. Again, the biases are represented as weights from an extra input  $x_0 = 1$ .

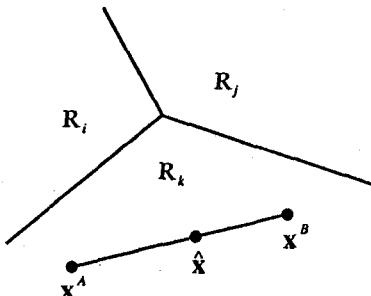


Figure 3.4. Example of decision boundaries produced by a multiclass linear discriminant. If two points  $\mathbf{x}^A$  and  $\mathbf{x}^B$  both lie in decision region  $\mathcal{R}_k$  then every point  $\hat{\mathbf{x}}$  on the line connecting them must also lie in region  $\mathcal{R}_k$ . It therefore follows that the decision regions must be simply connected and convex.

this, consider two points  $\mathbf{x}^A$  and  $\mathbf{x}^B$  which both lie in the region  $\mathcal{R}_k$  as shown in Figure 3.4. Any point  $\hat{\mathbf{x}}$  which lies on the line joining  $\mathbf{x}^A$  and  $\mathbf{x}^B$  can be written as

$$\hat{\mathbf{x}} = \alpha \mathbf{x}^A + (1 - \alpha) \mathbf{x}^B \quad (3.9)$$

where  $0 \leq \alpha \leq 1$ . Since  $\mathbf{x}^A$  and  $\mathbf{x}^B$  both lie in  $\mathcal{R}_k$ , they must satisfy  $y_k(\mathbf{x}^A) > y_j(\mathbf{x}^A)$  and  $y_k(\mathbf{x}^B) > y_j(\mathbf{x}^B)$  for all  $j \neq k$ . Using (3.4) and (3.9) it follows that  $y_k(\hat{\mathbf{x}}) = \alpha y_k(\mathbf{x}^A) + (1 - \alpha) y_k(\mathbf{x}^B)$  and hence  $y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}})$  for all  $j \neq k$ . Thus, all points on the line connecting  $\mathbf{x}^A$  and  $\mathbf{x}^B$  also lie in  $\mathcal{R}_k$  and so the region  $\mathcal{R}_k$  must be simply connected and convex.

$$\begin{aligned} \text{Circled } \alpha y_k(\mathbf{x}^A) &> \alpha y_j(\mathbf{x}^A) \quad \text{and} \quad (1-\alpha) y_k(\mathbf{x}^B) > (1-\alpha) y_j(\mathbf{x}^B) \quad \forall j \neq k \\ \Leftrightarrow \alpha y_k(\mathbf{x}^A) + (1-\alpha) y_k(\mathbf{x}^B) &> \alpha y_j(\mathbf{x}^A) + (1-\alpha) y_j(\mathbf{x}^B) \end{aligned}$$

### 3.1.3 Logistic discrimination

So far we have considered discriminant functions which are simple linear functions of the input variables. There are several ways in which such functions can be generalized, and here we consider the use of a non-linear function  $g(\cdot)$  which acts on the linear sum to give a discriminant function for the two-class problem of the form

$$y = g(\mathbf{w}^T \mathbf{x} + w_0) \quad (3.10)$$

where  $g(\cdot)$  is called an *activation function* and is generally chosen to be monotonic. The form (3.10) is still regarded as a linear discriminant since the decision boundary which it generates is still linear, as a consequence of the monotonic nature of  $g(\cdot)$ .

As a motivation for this form of discriminant, consider a two-class problem in which the class-conditional densities are given by Gaussian distributions with equal covariance matrices  $\Sigma_1 = \Sigma_2 = \Sigma$ , so that

$$p(\mathbf{x}|C_k) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}. \quad (3.11)$$

Using Bayes' theorem, the posterior probability of membership of class  $C_1$  is given by

$$P(C_1|\mathbf{x}) = \frac{p(\mathbf{x}|C_1)P(C_1)}{p(\mathbf{x}|C_1)P(C_1) + p(\mathbf{x}|C_2)P(C_2)} \quad (3.12)$$

$$= \frac{1}{1 + \exp(-a)} \quad (3.13)$$

$$= g(a) \quad (3.14)$$

where

$$a = \ln \frac{p(\mathbf{x}|C_1)P(C_1)}{p(\mathbf{x}|C_2)P(C_2)} \quad (3.15)$$

and the function  $g(a)$  is the *logistic sigmoid* activation function given by

$$g(a) \equiv \frac{1}{1 + \exp(-a)} \quad (3.16)$$

which is plotted in Figure 3.5. If we now substitute expressions for the class-conditional densities from (3.11) into (3.15) we obtain

*this is a lot of algebra, but looks straight forward*

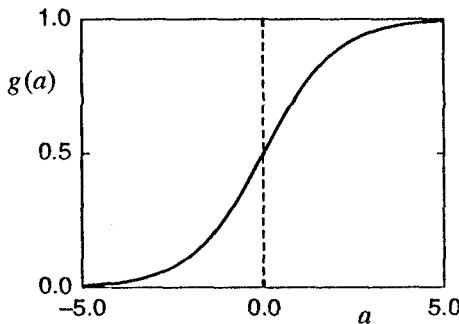


Figure 3.5. Plot of the logistic sigmoid activation function given by (3.16).

$$a = \mathbf{w}^T \mathbf{x} + w_0 \quad (3.17)$$

where

$$\mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2) \quad (3.18)$$

$$w_0 = -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{P(C_1)}{P(C_2)}. \quad (3.19)$$

Thus, we see that the use of the logistic sigmoid activation function allows the outputs of the discriminant to be interpreted as posterior probabilities. This implies that such a discriminant is providing more than simply a classification decision, and is potentially a very powerful result. The importance of interpreting the outputs of networks in terms of probabilities is discussed at much greater length in Chapter 6.

The term sigmoid means ‘S-shaped’, and the logistic form of the sigmoid maps the interval  $(-\infty, \infty)$  onto  $(0, 1)$ . If  $|a|$  is small, then the logistic sigmoid function  $g(a)$  can be approximated by a linear function, and so in this sense a network with sigmoidal activation functions contains a linear network as a special case. If there are more than two classes then an extension of the previous analysis leads to a generalization of the logistic sigmoid called a *normalized exponential* or *softmax*, which is discussed in detail in Section 6.9.

Linear discriminants with logistic activation functions have been widely used in the statistics literature under the name *logistic discrimination* (Anderson, 1982). Sigmoidal activation functions also play a crucial role in multi-layer neural networks, as discussed in Chapter 4.

Another form of linear discriminant was introduced by McCulloch and Pitts (1943) as a simple mathematical model for the behaviour of a single neuron in

a biological nervous system. Again this takes the form (3.10) with an activation function which is the **Heaviside step function**

$$g(a) = \begin{cases} 0 & \text{when } a < 0 \\ 1 & \text{when } a \geq 0. \end{cases} \quad (3.20)$$

In this model the inputs  $x_i$  represent the level of activity of other neurons which connect to the neuron being modelled, the weights  $w_i$  represent the strengths of the interconnections, called synapses, between the neurons, and the bias  $w_0$  represents the threshold for the neuron to 'fire'. Although this model has its origins in biology, it is clear that it can equally well be motivated within the framework of statistical pattern recognition. Networks of threshold units were studied by Rosenblatt (1962) under the name *perceptrons* and by Widrow and Hoff (1960) who called them *adalines*. They will be discussed in detail in Section 3.5.

Note that it is sometimes convenient to regard the linear discriminant (3.1) as a special case of the more general form (3.10). In this case the model is said to have a linear activation function, which in fact is just the identity  $g(a) \equiv a$ .

### 3.1.4 Binary input vectors

Linear discriminants, and the logistic activation function, also arise in a natural way when we consider input patterns in which the variables are binary (so that each  $x_i$  can take only the values 0 or 1). Let  $P_{ki}$  denote the probability that the input  $x_i$  takes the value +1 when the input vector is drawn from the class  $\mathcal{C}_k$ . The corresponding probability that  $x_i = 0$  is then given by  $1 - P_{ki}$ . We can combine these together to write the probability for  $x_i$  to take either of its allowed values in the form

$$p(x_i|\mathcal{C}_k) = P_{ki}^{x_i}(1 - P_{ki})^{1-x_i} \quad (3.21)$$

which is called a *Bernoulli* distribution. If we now assume that the input variables are statistically independent, we obtain the probability for the complete input vector as the product of the probabilities for each of the components separately:

$$p(\mathbf{x}|\mathcal{C}_k) = \prod_{i=1}^d P_{ki}^{x_i}(1 - P_{ki})^{1-x_i}. \quad (3.22)$$

We now recall from Chapter 1 that we can write a discriminant function which minimizes the probability of misclassifying new inputs in the form

$$y_k(\mathbf{x}) = \ln P(\mathbf{x}|\mathcal{C}_k) + \ln P(\mathcal{C}_k). \quad (3.23)$$

Substituting (3.22) into (3.23) we obtain a linear discriminant function given by

*Lots of algebra, but straightforward*

$$y_k(\mathbf{x}) = \sum_{i=1}^d w_{ki}x_i + w_{k0} \quad (3.24)$$

in which the weights and bias are given by

$$w_{ki} = \ln P_{ki} - \ln(1 - P_{ki}) \quad (3.25)$$

$$w_{k0} = \sum_{i=1}^d \ln(1 - P_{ki}) + \ln P(\mathcal{C}_k). \quad (3.26)$$

We have already seen that, for two classes with normally distributed class-conditional densities, the posterior probabilities can be obtained from the linear discriminant by applying a logistic activation function. A similar result holds also for the Bernoulli distribution. Consider a set of independent binary variables  $x_i$ , having Bernoulli class-conditional densities given by (3.22). If we substitute (3.22) into (3.12) we again obtain a single-layer network structure, with a logistic activation function, of the form

$$P(\mathcal{C}_1|\mathbf{x}) = g(\mathbf{w}^T\mathbf{x} + w_0) \quad (3.27)$$

where  $g(a)$  is given by (3.16) and

$$w_0 = \sum_i \ln \frac{1 - P_{1i}}{1 - P_{2i}} + \ln \frac{P(\mathcal{C}_1)}{P(\mathcal{C}_2)} \quad (3.28)$$

$$w_i = \ln \frac{P_{1i}}{P_{2i}} - \ln \frac{1 - P_{1i}}{1 - P_{2i}}. \quad (3.29)$$

We have shown that, both for normally distributed and Bernoulli distributed class-conditional densities, the posterior probabilities are obtained by a logistic single-layer network. In fact these are particular instances of a much more general result, which is derived in Section 6.7.1.

### 3.2 Linear separability

So far in this chapter we have discussed discriminant functions having a decision boundary which is linear, or more generally hyperplanar in higher dimensions. Clearly this is a very restricted class of decision boundary, and we might well expect such systems to have less than optimal performance for many practical applications. Indeed, this provides the principal motivation for using multi-layer networks of the kind discussed in Chapters 4 and 5. The particular nature of the limitation inherent in single-layer systems warrants some careful discussion, however.

in practice couldn't this be solved by adding another parameter  
 $x_3 = x_1 - x_2$

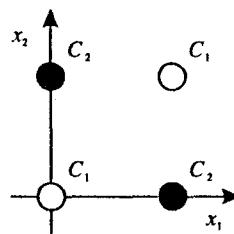


Figure 3.6. The exclusive-OR problem consists of four patterns in a two-dimensional space as shown. It provides a simple example of a problem which is not linearly separable.

Consider for the moment the problem of learning to classify a given data set exactly, where each input vector has been labelled as belonging to one of two classes  $C_1$  and  $C_2$ . If all of the points can be classified correctly by a linear (i.e. hyperplanar) decision boundary, then the points are said to be *linearly separable*. For such a data set there exist weight and bias values such that a linear discriminant will lead to perfect classification. A simple example of a data set which is not linearly separable is provided by the two-dimensional *exclusive-OR* problem, also known as *XOR*, illustrated in Figure 3.6. The input vectors  $\mathbf{x} = (0, 0)$  and  $(1, 1)$  belong to class  $C_1$ , while the input vectors  $(0, 1)$  and  $(1, 0)$  belong to class  $C_2$ . It is clear that there is no linear decision boundary which can classify all four points correctly. This problem can be generalized to  $d$ -dimensions when it is known as the  $d$ -bit parity problem. In this case the data set consists of all possible binary input vectors of length  $d$ , which are classified as class  $C_1$  if there is an even number of 1's in the input vector, and as class  $C_2$  otherwise.

For the case of continuous input variables it is interesting to consider the probability that a random set of patterns will be linearly separable. Suppose we have  $N$  data points distributed at random in  $d$  dimensions. Note that the particular distribution used to generate the random points is not relevant. All that we require is that there are no accidental degeneracies, i.e. that there is no subset of  $d$  or fewer points which are linearly dependent. The points are then said to be in *general position*. Having chosen the points, imagine that we then randomly assign each of the points to one of the two classes  $C_1$  and  $C_2$  with equal probability. Each possible assignment for the complete data set is referred to as a *dichotomy*, and for  $N$  points there are  $2^N$  possible dichotomies. We now ask what fraction  $F(N, d)$  of these dichotomies is linearly separable. It can be shown (Cover, 1965) that this fraction is given by the expression

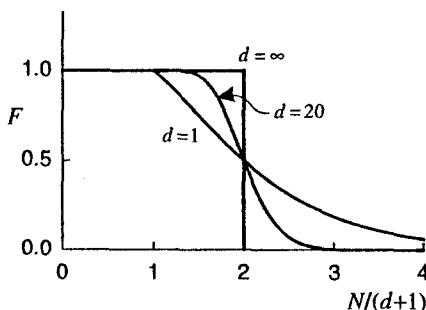


Figure 3.7. Plot of the fraction  $F(N, d)$  of the dichotomies of  $N$  data points in  $d$  dimensions which are linearly separable, as a function of  $N/(d + 1)$ , for various values of  $d$ .

$$F(N, d) = \begin{cases} 1 & \text{when } N \leq d + 1 \\ \frac{1}{2^{N-1}} \sum_{i=0}^d \binom{N-1}{i} & \text{when } N \geq d + 1 \end{cases} \quad (3.30)$$

which is plotted as a function of  $N/(d + 1)$  in Figure 3.7 for  $d = 1$ ,  $d = 20$  and  $d = \infty$ . Here the symbol

$$\binom{N}{M} \equiv \frac{N!}{(N - M)!M!} \quad (3.31)$$

denotes the number of combinations of  $M$  objects selected from a total of  $N$ . We see from (3.30) that, if the number of data points is fewer than  $d + 1$ , any labelling of the points will always lead to a linearly separable problem. For  $N = 2(d + 1)$ , the probability of linear separability is 0.5 for any value of  $d$  (Exercise 3.5). In a practical application, the positions of points from the same class will tend to be correlated, and so the probability that a data set with a much larger number of points than  $2(d + 1)$  will be linearly separable is higher than (3.30) would suggest.

For the case of binary input patterns, if there are  $d$  inputs then there are  $2^d$  possible input patterns and hence  $2^{2^d}$  possible labellings of those patterns between two classes. Those which can be implemented by a perceptron are called *threshold logic functions* and form an extremely small subset (less than  $2^{d^2}/d!$ ) of the total (Lewis and Coates, 1967).

In the neural computing literature a lot of attention is often paid to the inability of single-layer networks to solve simple problems such as XOR. From our statistical pattern recognition perspective, however, we see that the ability of a particular model to provide an exact representation of a given training set is

largely irrelevant. We are primarily interested in designing systems with good generalization performance, so that they give the greatest accuracy when presented with previously unseen data. Furthermore, problems such as XOR and parity involve learning the complete set of all possible input patterns, so the concept of generalization does not even apply. Finally, they have the property that the smallest possible change in the input pattern produces the largest possible change in the output. Most practical pattern recognition problems have the opposite characteristic, so that small changes in the inputs do not, for the most part, produce large changes in the outputs, and hence the mapping represented by the network should be relatively smooth.

Consider the problem of two normally-distributed classes with equal covariance matrices, discussed in Section 2.1.3. Since the class distributions overlap it is entirely possible that a finite sized data set drawn from these distributions will not be linearly separable. However, we know that the optimal decision boundary is in fact linear. A single-layer network can therefore achieve the best possible classification performance on unseen data, even though it may not separate the training data exactly.

The key consideration concerns the choice of an appropriate discriminant function for the particular problem in hand. This may involve a combination of prior knowledge of the general form which the solution should take, coupled with an empirical comparison of the performance of alternative models. These issues are considered in more detail in Chapters 8, 9 and 10. Here we simply note that single-layer networks correspond to a very narrow class of possible discriminant functions, and in many practical situations may not represent the optimal choice. Nevertheless, single-layer networks remain of considerable practical importance in providing a benchmark against which the performance of more complex multi-layer networks can be assessed. The fact that single-layer networks can often be trained very quickly, as shown in Section 3.4, gives them a particular advantage over more complex network structures which often require considerable computational effort to train.

### 3.3 Generalized linear discriminants

One way to generalize the discriminant functions, so as to permit a much larger range of possible decision boundaries, is to transform the input vector  $\mathbf{x}$  using a set of  $M$  predefined non-linear functions  $\phi_j(\mathbf{x})$ , sometimes called *basis functions*, and then to represent the output as a linear combination of these functions

$$y_k(\mathbf{x}) = \sum_{j=1}^M w_{kj} \phi_j(\mathbf{x}) + w_{k0}. \quad (3.32)$$

This now represents a much larger class of functions  $y_k(\mathbf{x})$ . In fact, as discussed in Chapters 4 and 5, for a suitable choice of the basis functions  $\phi_j(\mathbf{x})$ , the function in (3.32) can approximate any continuous functional transformation to arbitrary

accuracy. Again, we can absorb the biases as special cases of the weights by defining an extra basis function  $\phi_0 = 1$ , so that

$$y_k(\mathbf{x}) = \sum_{j=0}^M w_{kj} \phi_j(\mathbf{x}). \quad (3.33)$$

We have assumed that the basis functions  $\phi_j(\mathbf{x})$  are fixed, independently of the data. Chapters 4 and 5 discuss multi-layer neural networks, many of which can be regarded as generalized discriminant functions of the form (3.32), but in which the basis functions themselves can be modified during the training process.

### 3.4 Least-squares techniques

So far in this chapter we have discussed various forms of single-layer network and explored some of their properties. The remainder of the chapter is concerned with techniques for training such networks, and we begin with a discussion of methods based on the minimization of a sum-of-squares error function. This is the simplest form of error function and is most suitable for regression problems. While it can also be used for classification problems, there exist other, more appropriate, error functions, discussed at length in Chapter 6.

#### 3.4.1 Sum-of-squares error function

For consistency with the discussions in Chapter 5, we shall consider the error minimization problem in the context of the generalized linear network (3.33). This contains the simple linear discriminant of (3.4) as a special case in which the  $\phi_j(\mathbf{x})$  simply correspond to the input variables  $x_i$ . The sum-of-squares error function is given by a sum over all patterns in the training set, and over all outputs, of the form

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^c \{y_k(\mathbf{x}^n; \mathbf{w}) - t_k^n\}^2 \quad (3.34)$$

where  $y_k(\mathbf{x}^n; \mathbf{w})$  represents the output of unit  $k$  as a function of the input vector  $\mathbf{x}^n$  and the weight vector  $\mathbf{w}$ ,  $N$  is the number of training patterns, and  $c$  is the number of outputs. The quantity  $t_k^n$  represents the target value for output unit  $k$  when the input vector is  $\mathbf{x}^n$ . This error function is a smooth function of the weight parameters  $w_{kj}$ , and can be minimized by a variety of standard techniques. Since (3.33) is a linear function of the weights, the error function  $E(\mathbf{w})$  is a quadratic function of the weights, and hence its derivatives with respect to the weights are linear functions of the weights. The solution for the weight values at the minimum of the error function can therefore be found exactly in closed form, as we shall see in Section 3.4.3.

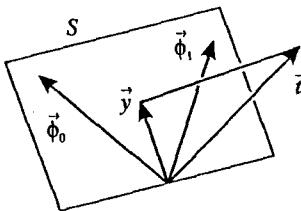


Figure 3.8. Geometrical interpretation of the solution to the least-squares problem, illustrated for the case of 3 training patterns ( $N = 3$ ) and 2 basis functions  $\phi_0$  and  $\phi_1$  (corresponding to  $M = 1$ ). The target values  $t^n$  are grouped together to form an  $N$ -dimensional vector  $\vec{t}$  which lives in an  $N$ -dimensional Euclidean space. The corresponding network outputs can similarly be represented as a vector  $\vec{y}$  which consists of a linear combination of  $M + 1$  basis vectors  $\vec{\phi}_j$ , which themselves span an  $(M + 1)$ -dimensional Euclidean sub-space  $S$ . The least-squares solution for  $\vec{y}$  is given by the orthogonal projection of  $\vec{t}$  onto  $S$ .

#### 3.4.2 Geometrical interpretation of least squares

Before deriving a solution for the weights, it is instructive to consider a geometrical interpretation of the least-squares problem. To do this we consider a network having a single output  $y$ . There is no loss of generality in doing this as the same discussion applies separately to each output of the network. For a particular input pattern  $x^n$  we can write the network output as

$$y^n = \sum_{j=0}^M w_j \phi_j^n \quad (3.35)$$

where  $\phi_j^n \equiv \phi_j(x^n)$ . We now group the target values together to form an  $N$ -dimensional vector  $\vec{t}$  whose elements are given by  $t^n$ . This vector can be considered to live in an  $N$ -dimensional Euclidean space, as indicated in Figure 3.8. For each basis function  $\phi_j(x)$  we can similarly group the  $N$  values of  $\phi_j^n$ , corresponding to the  $N$  data points, to make a vector  $\vec{\phi}_j$ , also of dimension  $N$ , which can be drawn in the same space as the vector  $\vec{t}$ . For the moment we shall assume that the number of basis functions (including the bias) is less than the number of patterns, so that  $M + 1 < N$ . The  $M + 1$  vectors  $\vec{\phi}_j$ , corresponding to the  $M + 1$  basis functions, then form a (non-orthogonal) basis set which spans an  $(M + 1)$ -dimensional Euclidean sub-space  $S$ . The network outputs  $y^n$  can also be grouped to form a vector  $\vec{y}$ . From (3.35) we see that  $\vec{y}$  is given by a linear combination of the  $\vec{\phi}_j$  of the form

$$\vec{y} = \sum_{j=0}^M w_j \vec{\phi}_j \quad (3.36)$$

so that  $\vec{y}$  is constrained to lie in the sub-space  $S$ , as shown in Figure 3.8. By changing the values of the weights  $w_j$  we can change the location of  $\vec{y}$  subject to this constraint.

The sum-of-squares error (3.34) can now be written in the form

$$E = \frac{1}{2} \left\| \sum_{j=0}^M w_j \vec{\phi}_j - \vec{t} \right\|^2 \quad (3.37)$$

If we minimize this expression with respect to the weights  $w_j$  we find

$$\frac{\partial E}{\partial w_j} = 0 = \vec{\phi}_j^T (\vec{y} - \vec{t}), \quad j = 1, \dots, M. \quad (3.38)$$

This represents a set of coupled equations for the weights, known as the *normal equations* of the least-squares problem, for which we shall find an explicit solution shortly. Before doing so, however, it is useful to consider the geometrical interpretation of (3.38). Let us decompose  $\vec{t}$  into the sum of two vectors  $\vec{t} = \vec{t}_{\perp} + \vec{t}_{\parallel}$  where  $\vec{t}_{\parallel}$  is the orthogonal projection of  $\vec{t}$  onto the sub-space  $S$ , and  $\vec{t}_{\perp}$  is the remainder. Then  $\vec{\phi}_j^T \vec{t}_{\perp} = 0$  by definition, and hence from (3.38) we have

$$\vec{\phi}_j^T (\vec{y} - \vec{t}_{\parallel}) = 0, \quad j = 1, \dots, M. \quad (3.39)$$

Since the vectors  $\vec{\phi}_j$  form a basis set which span the sub-space  $S$ , we can solve (3.39) to give

$$\vec{y} = \vec{t}_{\parallel} \quad (3.40)$$

**and so the solution vector is just the projection of the vector of target values onto the sub-space spanned by the basis vectors, as indicated in Figure 3.8.** This result is intuitively correct, since the process of learning corresponds to choosing a direction for  $\vec{y}$  such as to minimize its distance from  $\vec{t}$ . Since  $\vec{y}$  is constrained to lie in the sub-space, the best we can do is choose it to correspond to the orthogonal projection of  $\vec{t}$  onto  $S$ . This minimizes the length of the error vector  $\vec{\epsilon} = \vec{y} - \vec{t}$ . Note that the residual error vector  $\vec{\epsilon}_{\min} = \vec{t}_{\parallel} - \vec{t} = -\vec{t}_{\perp}$  is then orthogonal to  $S$ , so that  $\vec{\phi}_j^T \vec{\epsilon}_{\min} = 0$ .

### 3.4.3 Pseudo-inverse solution

We now proceed to find an exact solution to the least-squares problem. To do this we return to the case of a network having  $c$  outputs. Using the expression (3.33), we can write the sum-of-squares error function (3.34) in the form

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^c \left\{ \sum_{j=0}^M w_{kj} \phi_j^n - t_k^n \right\}^2. \quad (3.41)$$

Differentiating this expression with respect to  $w_{kj}$  and setting the derivative to zero gives the normal equations for the least-squares problem in the form

$$\sum_{n=1}^N \left\{ \sum_{j'=0}^M w_{kj'} \phi_{j'}^n - t_k^n \right\} \phi_j^n = 0. \quad (3.42)$$

In order to find a solution to (3.42) it is convenient to write it in a matrix notation to give

$$(\Phi^T \Phi) \mathbf{W}^T = \Phi^T \mathbf{T}. \quad (3.43)$$

Here  $\Phi$  has dimensions  $N \times M$  and elements  $\phi_j^n$ ,  $\mathbf{W}$  has dimensions  $c \times M$  and elements  $w_{kj}$ , and  $\mathbf{T}$  has dimensions  $N \times c$  and elements  $t_k^n$ . The matrix  $\Phi^T \Phi$  in (3.43) is a square matrix of dimension  $M \times M$ . Provided it is non-singular we may invert it to obtain a solution to (3.43) which can be written in the form

$$\mathbf{W}^T = \Phi^\dagger \mathbf{T} \quad (3.44)$$

where  $\Phi^\dagger$  is an  $M \times N$  matrix known as the *pseudo-inverse* of  $\Phi$  (Golub and Kahan, 1965; Rao and Mitra, 1971) and is given by

$$\Phi^\dagger \equiv (\Phi^T \Phi)^{-1} \Phi^T \quad (3.45)$$

Since  $\Phi$  is, in general, a non-square matrix it does not itself have a true inverse, but the pseudo-inverse does have the property (as is easily seen from 3.45) that  $\Phi^\dagger \Phi = \mathbf{I}$  where  $\mathbf{I}$  is the unit matrix. Note, however, that  $\Phi \Phi^\dagger \neq \mathbf{I}$  in general. If the matrix  $\Phi^T \Phi$  is singular then (3.43) does not have a unique solution. However, if the pseudo-inverse is defined by

$$\Phi^\dagger \equiv \lim_{\epsilon \rightarrow 0} (\Phi^T \Phi + \epsilon \mathbf{I})^{-1} \Phi^T \quad (3.46)$$

then it can be shown that the limit always exists, and that this limiting value minimizes  $E$  (Rao and Mitra, 1971).

In practice, the direct solution of the normal equations can lead to numerical difficulties due to the possibility of  $\Phi^T \Phi$  being singular or nearly singular. This can arise if two of the basis vectors  $\vec{\phi}_j$ , shown in Figure 3.8, are nearly collinear. The effects of noise and numerical error can then lead to very large values for the weights which give near cancellation between these vectors. Figure 3.9(a) shows two basis vectors  $\vec{\phi}_1$  and  $\vec{\phi}_2$  which are nearly orthogonal, together with the component  $\vec{y}_{\parallel}$  of  $\vec{y}$  which lies in the plane spanned by  $\vec{\phi}_1$  and  $\vec{\phi}_2$ . The corresponding weight values needed to express  $\vec{y}_{\parallel}$  as a linear combination of  $\vec{\phi}_1$  and  $\vec{\phi}_2$  have relatively small values. By contrast, Figure 3.9(b) shows the corresponding situation when the vectors  $\vec{\phi}_1$  and  $\vec{\phi}_2$  are nearly collinear. In this case the weights need to adopt large (positive or negative) values in order to represent  $\vec{y}_{\parallel}$  as a linear combination of the basis vectors. In the case where the two basis vectors are exactly collinear, we can write  $\vec{\phi}_2 = \lambda \vec{\phi}_1$  for some constant  $\lambda$ . Then  $w_1 \vec{\phi}_1 + w_2 \vec{\phi}_2 = (w_1 + \lambda w_2) \vec{\phi}_1$  and only the combination  $(w_1 + \lambda w_2)$  is fixed by the least-squares procedure, with the value of  $w_2$ , say, being arbitrary. Near degeneracies will not be uncommon when dealing with real, noisy data sets. In practice, such problems are best resolved by using the technique of *singular value decomposition* (SVD) to find a solution for the weights. A good introduction to SVD, together with a suggested numerical implementation, can be found in Press *et al.* (1992). Such an approach avoids problems due to the accumulation of numerical roundoff errors, and automatically selects (from amongst a set of nearly degenerate solutions) the one for which the length  $\|w_k\|$  of the  $k$ th weight vector is shortest.

In the above discussion, the bias parameters were treated as a special case of the weights. We can gain some insight into the role of the biases if we make them explicit. If we consider the minimization of (3.41) with respect to the bias parameters alone we obtain

$$\frac{\partial E}{\partial w_{k0}} = \sum_{n=1}^N \left\{ \sum_{j=1}^M w_{kj} \phi_j^n + w_{k0} - t_k^n \right\} = 0 \quad (3.47)$$

which can be solved for the biases to give

$$w_{k0} = \bar{t}_k - \sum_{j=1}^M w_{kj} \bar{\phi}_j \quad (3.48)$$

where

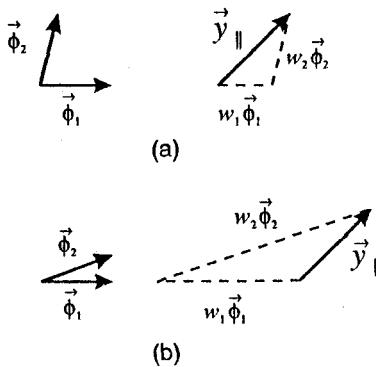


Figure 3.9. In (a) we see two basis vectors  $\vec{\phi}_1$  and  $\vec{\phi}_2$  which are nearly orthogonal. The least-squares solution vector  $\vec{y}_{\parallel}$  is given by a linear combination of these vectors, with relatively small values for the coefficients  $w_1$  and  $w_2$ . In (b) the basis vectors are nearly collinear, and the magnitudes of the corresponding weight values become very large.

$$\bar{t}_k = \frac{1}{N} \sum_{n=1}^N t_k^n, \quad \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j^n. \quad (3.49)$$

This result tells us that the role of the bias parameters is to compensate for the difference between the mean (over the training set) of the output vector for the network and the corresponding mean of the target data.

If  $\Phi^T$  is a square non-singular matrix, the pseudo-inverse reduces to the usual inverse. The matrix is square when  $N = M$ , so that the number of patterns equals the number of basis functions. If we multiply (3.43) by  $(\Phi^T)^{-1}$  we obtain

$$\Phi \mathbf{W}^T = \mathbf{T}. \quad (3.50)$$

If we write this in index notation we have

$$\sum_{j=0}^M w_{kj} \phi_j^n = t_k^n \quad (3.51)$$

and we see that, for each input pattern, the network outputs are exactly equal to the corresponding target values, and hence the sum-of-squares error (3.41) will be zero. The condition for  $(\Phi^T)^{-1}$  to exist is that the columns  $\phi^n$  of the matrix  $\Phi^T$  be linearly independent. If the vectors  $\phi^n$  are not linearly independent, so that the effective value of  $N$  is less than  $M$ , then the least-squares problem

is under-determined. Similarly, if there are fewer patterns than basis functions, so that  $N < M$ , then the least-squares problem is again under-determined. In such cases, there is a continuum of solutions for the weights, all of which give zero error. Singular value decomposition leads to a numerically well-behaved algorithm which picks out the particular solution for which the magnitude  $\|\mathbf{w}_k\|$  of the weight vector for each output unit  $k$  is the shortest. As we have already indicated in Chapter 1, it is desirable to have a sufficiently large training set that the weight values are ‘over-determined’, so that in practice we arrange that  $N > M$ , which corresponds to the situation depicted in Figure 3.8.

### 3.4.4 Gradient descent

We have shown how, for a linear network, the weight values which minimize the sum-of-squares error function can be found explicitly in terms of the pseudo-inverse of a matrix. It is important to note that this result is only possible for the case of a linear network, with a sum-of-squares error function. If a non-linear activation function, such as a sigmoid, is used, or if a different error function is considered, then a closed form solution is no longer possible. However, if the activation function is differentiable, as is the case for the logistic sigmoid in (3.16) for instance, the derivatives of the error function with respect to the weight parameters can easily be evaluated. These derivatives can then be used in a variety of gradient-based optimization algorithms, discussed in Chapter 7, for finding the minimum of the error function. Here we consider one of the simplest of such algorithms, known as gradient descent.

It is convenient to group all of the parameters (weights and biases) in the network together to form a single weight vector  $\mathbf{w}$ , so that the error function can be expressed as  $E = E(\mathbf{w})$ . Provided  $E$  is a differentiable function of  $\mathbf{w}$  we may adopt the following procedure. We begin with an initial guess for  $\mathbf{w}$  (which might for instance be chosen at random) and we then update the weight vector by moving a small distance in  $\mathbf{w}$ -space in the direction in which  $E$  decreases most rapidly, i.e. in the direction of  $-\nabla_{\mathbf{w}}E$ . By iterating this process we generate a sequence of weight vectors  $\mathbf{w}^{(\tau)}$  whose components are calculated using

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}} \quad (3.52)$$

where  $\eta$  is a small positive number called the *learning rate* parameter. Under suitable conditions the sequence of weight vectors will converge to a point at which  $E$  is minimized. The choice of the value for  $\eta$  can be fairly critical, since if it is too small the reduction in error will be very slow, while, if it is too large, divergent oscillations can result.

In general the error function is given by a sum of terms each of which is calculated using just one of the patterns from the training set, so that

$$E(\mathbf{w}) = \sum_n E^n(\mathbf{w}) \quad (3.53)$$

where the term  $E^n$  is calculated using pattern  $n$  only. In this case we can update the weight vector using just one pattern at a time

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \frac{\partial E^n}{\partial w_{kj}} \quad (3.54)$$

and this is repeated many times by cycling through all of the patterns used in the definition of  $E$ . This form of sequential, or pattern-based, update is reminiscent of the Robbins–Monro procedure introduced in Section 2.4, and many of the same comments apply here. In particular, this technique allows the system to be used in real-time adaptive applications in which data is arriving continuously. Each data point can be used once and then discarded, and if the value of  $\eta$  is chosen appropriately, the system may be able to ‘track’ any slow changes in the characteristics of the data. If  $\eta$  is chosen to decrease with time in a suitable way during the learning process, then gradient descent becomes precisely the Robbins–Monro procedure for finding the root of the regression function  $\mathcal{E}[\partial E^n / \partial w_l]$  where  $\mathcal{E}$  denotes the expectation. If the value of  $\eta$  is chosen to be steadily decreasing with time, so that  $\eta^{(\tau)} = \eta_0 / \tau$  (which satisfies the conditions for the Robbins–Monro theorem stated in Section 2.4), then the weight matrix  $\mathbf{W}$  can be shown to converge to a solution of

$$\Phi^T (\Phi \mathbf{W} - \mathbf{T}) = 0 \quad (3.55)$$

where  $\Phi$  is defined on page 92, irrespective of whether or not  $\Phi^T \Phi$  is singular. Gradient descent, and its limitations, are discussed at greater length in Chapter 7, along with a variety of more sophisticated optimization algorithms.

In order to implement gradient descent, we need explicit expressions for the derivatives of the error function with respect to the weights. We consider first the pattern-based form of gradient descent given by (3.54). For a generalized linear network function of the form (3.33) the derivatives are given by

$$\frac{\partial E^n}{\partial w_{kj}} = \{y_k(\mathbf{x}^n) - t_k^n\} \phi_j(\mathbf{x}^n) = \delta_k^n \phi_j^n \quad (3.56)$$

where we have defined

$$\delta_k^n \equiv y_k(\mathbf{x}^n) - t_k^n. \quad (3.57)$$

We see that the derivative with respect to a weight  $w_{kj}$  connecting basis function  $j$  to output  $k$  can be expressed as the product of  $\delta_k$  for the output unit and  $\phi_j$  for the basis function. Thus, the derivative can be calculated from quantities which are ‘local’ (in the sense of the network diagram) to the weight concerned. This property is discussed at greater length in the context of multi-layer networks in Section 4.8. Combining (3.54) and (3.56) we see that the change in the weights

due to presentation of a particular pattern is given by

$$\Delta w_{kj} = -\eta \delta_k^n \phi_j^n. \quad (3.58)$$

This rule, and its variants, are known by a variety of names including the LMS (least mean squares) rule, the adaline rule, the Widrow–Hoff rule (Widrow and Hoff, 1960), and the delta rule.

For networks with differentiable non-linear activation functions, such as the logistic sigmoid shown in Figure 3.5, we can write the network outputs in the form

$$y_k = g(a_k) \quad (3.59)$$

where  $g(\cdot)$  is the activation function, and

$$a_k = \sum_{j=0}^M w_{kj} \phi_j. \quad (3.60)$$

The derivatives of the error function for pattern  $n$  again take the form

$$\frac{\partial E^n}{\partial w_{kj}} = g'(a_k) \delta_k^n \phi_j^n \quad (3.61)$$

in which

$$\delta_k^n = g'(a_k)(y_k(x^n) - t_k^n). \quad (3.62)$$

For the logistic sigmoid given by (3.16), the derivative of the activation function can be expressed in the simple form

$$g'(a) = g(a)(1 - g(a)). \quad (3.63)$$

For gradient descent based on the total error function (summed over all patterns in the training set) given by (3.52), the derivatives are obtained by computing the derivatives for each pattern separately and then summing over all patterns

$$\frac{\partial E}{\partial w_{kj}} = \sum_n \frac{\partial E^n}{\partial w_{kj}}. \quad (3.64)$$

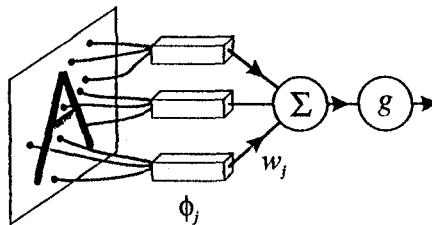


Figure 3.10. The perceptron network used a fixed set of processing elements, denoted  $\phi_j$ , followed by a layer of adaptive weights  $w_j$  and a threshold activation function  $g(\cdot)$ . The processing elements  $\phi_j$  typically also had threshold activation functions, and took inputs from a randomly chosen subset of the pixels of the input image.

### 3.5 The perceptron

Single-layer networks, with threshold activation functions, were studied by Rosenblatt (1962) who called them *perceptrons*. Rosenblatt also built hardware implementations of these networks, which incorporated learning using an algorithm to be discussed below. These networks were applied to classification problems, in which the inputs were usually binary images of characters or simple shapes. The properties of perceptrons are reviewed in Block (1962).

At the same time as Rosenblatt was developing the perceptron, Widrow and co-workers were working along similar lines using systems known as *adalines* (Widrow and Lehr, 1990). The term adaline comes from ADaptive LINear Element, and refers to a single processing unit with threshold non-linearity (Widrow and Hoff, 1960) of essentially the same form as the perceptron.

We have already seen that a network with a single layer of weights has very limited capabilities. To improve the performance of the perceptron, Rosenblatt used a layer of fixed processing elements to transform the raw input data, as shown in Figure 3.10. These processing elements can be regarded as the basis functions of a generalized linear discriminant. They typically took the form of fixed weights connected to a random subset of the input pixels, with a threshold activation function of the form (3.20). We shall again use the convention introduced earlier of defining an extra basis function  $\phi_0$  whose activation is permanently set to +1, together with a corresponding bias parameter  $w_0$ . The output of the perceptron is therefore given by

$$y = g \left( \sum_{j=0}^M w_j \phi_j(x) \right) = g(\mathbf{w}^T \boldsymbol{\phi}) \quad (3.65)$$

where  $\boldsymbol{\phi}$  denotes the vector formed from the activations  $\phi_0, \dots, \phi_M$ . The output

unit activation function is most conveniently chosen to be an anti-symmetric version of the threshold activation function of the form

$$g(a) = \begin{cases} -1 & \text{when } a < 0 \\ +1 & \text{when } a \geq 0. \end{cases} \quad (3.66)$$

We now turn to a discussion of the procedures used to train the perceptron.

### 3.5.1 The perceptron criterion

Since our goal is to produce an effective classification system, it would be natural to define the error function in terms of the total number of misclassifications over the training set. More generally we could introduce a loss matrix (Section 1.10) and consider the total loss incurred as a result of a particular classification of the data set. Such error measures, however, prove very difficult to work with in practice. This is because smooth changes in the values of the weights (and biases) cause the decision boundaries to move across the data points resulting in discontinuous changes in the error. The error function is therefore piecewise constant, and so procedures akin to gradient descent cannot be applied. We therefore seek other error functions which can be more easily minimized.

In this section we consider a continuous, piecewise-linear error function called the *perceptron criterion*. As each input vector  $x^n$  is presented to the inputs of the network it generates a corresponding vector of activations  $\phi^n$  in the first-layer processing elements. Suppose we associate with each input vector  $x^n$  a corresponding target value  $t^n$ , such that the desired output from the network is  $t^n = +1$  if the input vector belongs to class  $C_1$ , and  $t^n = -1$  if the vector belongs to class  $C_2$ . From (3.65) and (3.66) we want  $w^T \phi^n > 0$  for vectors from class  $C_1$ , and  $w^T \phi^n < 0$  for vectors from class  $C_2$ . It therefore follows that for all vectors we want to have  $w^T(\phi^n t^n) > 0$ . This suggests that we try to minimize the following error function, known as the perceptron criterion

$$E^{\text{perc}}(w) = - \sum_{\phi^n \in \mathcal{M}} w^T(\phi^n t^n) \quad (3.67)$$

where  $\mathcal{M}$  is the set of vectors  $\phi^n$  which are *misclassified* by the current weight vector  $w$ . The error function  $E^{\text{perc}}(w)$  is the sum of a number of positive terms, and equals zero if all of the data points are correctly classified. From the discussion in Section 3.1 we see that  $E^{\text{perc}}(w)$  is proportional to the sum, over all of the input patterns which are misclassified, of the (absolute) distances to the decision boundary. During training, the decision boundary will move and some points which were previously misclassified will become correctly classified (and vice versa) so that the set of patterns which contribute to the sum in (3.67) will change. The perceptron criterion is therefore continuous and piecewise linear with discontinuities in its gradient.

 doesn't the sum make the function itself discontinuous?

### 3.5.2 Perceptron learning

If we apply the pattern-by-pattern gradient descent rule (3.54) to the perceptron criterion (3.67) we obtain

$$w_j^{(\tau+1)} = w_j^{(\tau)} + \eta \phi_j^n t^n. \quad (3.68)$$

This corresponds to a very simple learning algorithm which can be summarized as follows. Cycle through all of the patterns in the training set and test each pattern in turn using the current set of weight values. If the pattern is correctly classified do nothing, otherwise add the pattern vector (multiplied by  $\eta$ ) to the weight vector if the pattern is labelled class  $C_1$  or subtract the pattern vector (multiplied by  $\eta$ ) from the weight vector if the pattern is labelled class  $C_2$ . It is easy to see that this procedure tends to reduce the error since

$$-\mathbf{w}^{(\tau+1)\top} (\phi^n t^n) = -\mathbf{w}^{(\tau)\top} (\phi^n t^n) - \eta (\phi^n t^n)^\top (\phi^n t^n) < -\mathbf{w}^{(\tau)\top} (\phi^n t^n) \quad (3.69)$$

*as long as we normalize initial parameters*

since  $\|\phi^n t^n\|^2 > 0$  and  $\eta > 0$ .

For the particular case of the perceptron criterion, we see that the value of  $\eta$  is in fact unimportant since a change in  $\eta$  is equivalent to a re-scaling of the weights and bias (assuming the initial parameter values are similarly re-scaled). This leaves the location of the decision boundaries unchanged. To see this, recall that the location of the decision boundary is given by (3.2), and is therefore unchanged if all of the weights, including the bias, are rescaled by the same constant. Thus, when minimizing the perceptron criterion, we can take  $\eta = 1$  with no loss of generality. This property does not hold, however, for most other forms of error function.

In Figures 3.11–3.13 we give a simple example of learning in a perceptron, for the case of one basis function  $\phi_1$ , so that, with biases included as special cases of the weights, the data points live in a two-dimensional space  $(\phi_0, \phi_1)$  with  $\phi_0 = 1$ .

### 3.5.3 Perceptron convergence theorem

There is an interesting result which states that, for any data set which is linearly separable, the learning rule in (3.68) is guaranteed to find a solution in a finite number of steps (Rosenblatt, 1962; Block, 1962; Nilsson, 1965; Minsky and Papert, 1969; Duda and Hart, 1973; Hand, 1981; Arbib, 1987; Hertz *et al.*, 1991). This is known as the *perceptron convergence theorem*. Here we give a relatively simple proof, based on Hertz *et al.* (1991).

Since we are considering a data set which is linearly separable, we know that there exists at least one weight vector  $\hat{\mathbf{w}}$  for which all training vectors are correctly classified, so that

$$\hat{\mathbf{w}}^\top \phi^n t^n > 0 \quad \text{for all } n. \quad (3.70)$$

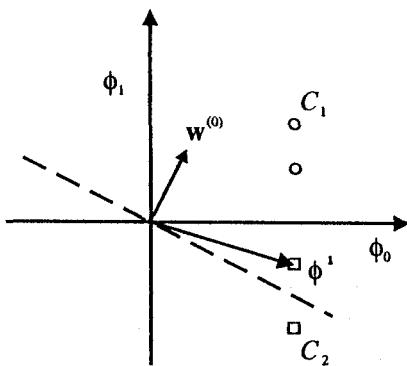


Figure 3.11. A simple example of perceptron learning, for a data set with four patterns. Circles represent patterns belonging to class  $C_1$  and squares represent patterns belonging to class  $C_2$ . The initial decision boundary, corresponding to the weight vector  $w^{(0)}$ , shown by the dashed curve, leaves one of the points, at  $\phi^1$ , incorrectly classified.

The learning process starts with some arbitrary weight vector which, without loss of generality, we can assume to be the zero vector. At each step of the algorithm, the weight vector is updated using

$$w^{(\tau+1)} = w^{(\tau)} + \phi^n t^n \quad (3.71)$$

where  $\phi^n$  is a vector which is misclassified by the perceptron. Suppose that, after running the algorithm for some time, the number of times that each vector  $\phi^n$  has been presented and misclassified is  $\tau^n$ . Then the weight vector at this point will be given by

$$w = \sum_n \tau^n \phi^n t^n. \quad (3.72)$$

We now take the scalar product of this equation with  $\hat{w}$  to give

$$\begin{aligned} \hat{w}^T w &= \sum_n \tau^n \hat{w}^T \phi^n t^n \\ &\geq \tau \min_n (\hat{w}^T \phi^n t^n) \end{aligned} \quad (3.73)$$

where  $\tau = \sum_n \tau^n$  is the total number of weight updates, and the inequality results from replacing each update vector by the smallest of the update vectors.

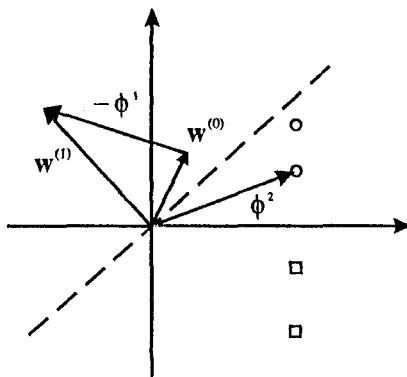


Figure 3.12. To correct for the misclassification of  $\phi^1$  in Figure 3.11 we add (minus)  $\phi^1$  onto  $w^{(0)}$  to give a new weight vector  $w^{(1)}$ , with the new decision boundary again shown by the dashed curve. The point at  $\phi^1$  is now correctly classified, but the point at  $\phi^2$  is now incorrectly classified.

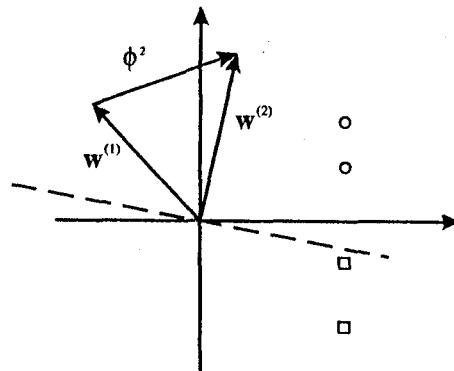


Figure 3.13. To correct for the misclassification of  $\phi^2$  in Figure 3.12 we add  $\phi^2$  onto  $w^{(1)}$  to give a new weight vector  $w^{(2)}$  which classifies all the points correctly.

From (3.70) it then follows that the value of  $\hat{w}^T w$  is bounded below by a function which grows linearly with  $\tau$ .

Keeping this result in mind, we now turn to a consideration of the magnitude of the weight vector  $w$ . From (3.71) we have

$$\begin{aligned}\|\mathbf{w}^{(\tau+1)}\|^2 &= \|\mathbf{w}^{(\tau)}\|^2 + \|\phi^n\|^2(t^n)^2 + 2\mathbf{w}^{(\tau)\top}\phi^n t^n \\ &\leq \|\mathbf{w}^{(\tau)}\|^2 + \|\phi^n\|^2(t^n)^2\end{aligned}\quad (3.74)$$

where the inequality follows from the fact that the pattern  $\phi^n$  must have been misclassified, and so  $\mathbf{w}^{(\tau)\top}\phi^n t^n < 0$ . We also have  $(t^n)^2 = 1$  since  $t^n = \pm 1$ , and  $\|\phi^n\|^2 \leq \|\phi\|_{\max}^2$  where  $\|\phi\|_{\max}$  is the length of the longest input vector. Thus, the change in the value of  $\|\mathbf{w}\|^2$  satisfies

$$\Delta \|\mathbf{w}\|^2 \equiv \|\mathbf{w}^{(\tau+1)}\|^2 - \|\mathbf{w}^{(\tau)}\|^2 \leq \|\phi\|_{\max}^2 \quad (3.75)$$

and so after  $\tau$  weight vector updates we have

$$\|\mathbf{w}\|^2 \leq \tau \|\phi\|_{\max}^2 \quad (3.76)$$

and so the length  $\|\mathbf{w}\|$  of the weight vector increases no faster than  $\tau^{1/2}$ . We now recall the previous result that  $\hat{\mathbf{w}}^\top \mathbf{w}$  is bounded below by a linear function of  $\tau$ . Since  $\hat{\mathbf{w}}$  is fixed, we see that for sufficiently large  $\tau$  these two results would become incompatible. Thus  $\tau$  cannot grow indefinitely, and so the algorithm must converge in a finite number of steps.

One of the difficulties with the perceptron learning rule is that, if the data set happens not to be linearly separable, then the learning algorithm will never terminate. Furthermore, if we arbitrarily stop the learning process there is no guarantee that the weight vector found will generalize well for new data. Various heuristics have been proposed with a view to giving good performance on problems which are not linearly separable while still ensuring convergence when the problem is linearly separable. For example, the value of the parameter  $\eta$  may be made to decrease during the learning process so that the corrections gradually become smaller. One approach is to take  $\eta = K/\tau$  where  $K$  is a constant and  $\tau$  is the step number, by analogy with the Robbins–Monro procedure (Section 2.4.1). An alternative algorithm for finding good solutions on problems which are not linearly separable, called the *pocket algorithm*, is described in Section 9.5.1. As we have already discussed, the issue of linear separability is a somewhat artificial one, and it is more important to develop learning algorithms which can be expected to give good performance across a wide range of problems, even if this means sacrificing the guarantee of perfect classification for linearly separable problems.

### 3.5.4 Limitations of the perceptron

When perceptrons were being studied experimentally in the 1960s, it was found that they could solve many problems very readily, whereas other problems, which superficially appeared to be no more difficult, proved impossible to solve. A critical appraisal of the capabilities of these networks, from a formal mathematical viewpoint, was given by Minsky and Papert (1969) in their book *Perceptrons*.

They showed that there are many types of problem which a perceptron cannot, in any practical sense, be used to solve. In this context a solution is taken to be a correct classification of all of the patterns in the training set.

Many recent textbooks on neural networks have summarized Minsky and Papert's contribution by pointing out that a single-layer network can only classify data sets which are linearly separable, and hence can not solve problems such as the XOR example considered earlier. In fact, the arguments of Minsky and Papert are rather more subtle, and shed light on the nature of multi-layer networks in which only one of the layers of weights is adaptive. Consider the perceptron shown in Figure 3.10. The first layer of fixed (non-adaptive) processing units computes a set of functions  $\phi_j$  whose values depend on the input pattern. Even though the data set of input patterns may not be linearly separable, when viewed in the space of original input variables, it can easily be the case that the same set of patterns becomes linearly separable when transformed into the space of  $\phi_j$  values. Thus a perceptron can solve a linearly inseparable problem, provided it has an appropriate set of first-layer processing elements.

The real difficulty with the perceptron arises from the fact that these processing elements are fixed in advance and cannot be adapted to the particular problem (or data set) which is being considered. As a consequence of this, it turns out that the number, or complexity, of such units must grow very rapidly (typically exponentially) with the dimensionality of the problem if the perceptron is to remain capable in general of providing a solution. It is therefore necessary to limit either the number or the complexity of the first-layer units. Minsky and Papert discuss a range of different forms of perceptron (depending on the form of the functions  $\phi_j$ ) and for each of them they provide examples of problems which cannot be solved.

Here we consider one particular form, called a *diameter-limited perceptron*, in which we consider two-dimensional input images as shown in Figure 3.10, and in which each of the  $\phi_j$  takes its inputs only from within a small localized region of the image, called a *receptive field*, having fixed diameter. Minsky and Papert (1969) provide a simple geometrical proof that such a perceptron cannot solve a simple problem involving the determination of whether a binary geometrical image is simply connected. This is illustrated in Figure 3.14. We shall suppose that connected shapes are labelled with targets +1 and that disconnected shapes have targets -1. Note that the overall length of the shapes is taken to be much larger than the maximum diameter of the receptive fields (indicated by the dashed circles), so that no single receptive field can overlap both ends of the shape. For the shape in Figure 3.14 (a), the functions  $\phi_j$  and the adaptive weights in the perceptron must be such that the linear sum which forms the input to the threshold function is negative, if this figure is to be correctly classified as 'disconnected'. In going to 3.14 (b), only the left-hand end of the shape has changed, so the receptive fields which lie in this region, and their corresponding weights, must be such that the linear sum is increased sufficiently to make it go positive, since this shape is 'connected'. Similarly, in going from 3.14 (a) to 3.14 (c) the linear sum must also be increased sufficiently to make it positive. However, in going

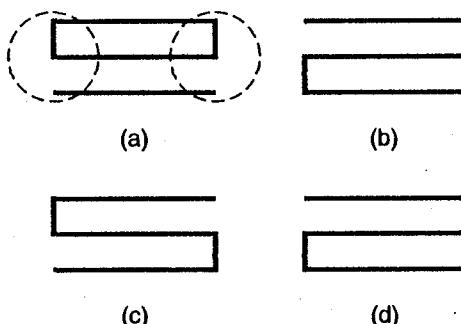


Figure 3.14. An example of a simple problem, involving the determination of whether a geometrical figure is simply connected, which cannot be solved by a perceptron whose inputs are taken from regions of limited diameter.

from 3.14 (a) to 3.14 (d), both ends of the shape have been changed in this way, and so the linear sum must be even more positive. This is inevitable since the diameter limitation means that the response due to the two ends of the shape are independent. Thus, the linear sum cannot be negative for the shape in 3.14 (d), which will therefore be misclassified.

Various alternative approaches to limiting the complexity of the first-layer units can be considered. For instance, in an *order-limited* perceptron, each of the  $\phi_j$  can take inputs only from a limited number of input pixels (which may lie anywhere on the input image). Counter-examples similar to the one presented above can be found also for these other choices of  $\phi_j$ . These difficulties can be circumvented by allowing the number and complexity of the  $\phi_j$  to grow sufficiently rapidly with the dimensionality of the problem. For example, it is shown in Section 4.2.1 that, for networks with binary inputs, there is a simple procedure for constructing the  $\phi_j$  such that any set of input patterns is guaranteed to be linearly separable in the  $\phi_j$  space. The number of such units, however, must grow exponentially with the input dimensionality. Such an approach is therefore totally impractical for anything other than toy problems.

The practical solution to these difficulties is to allow the functions  $\phi_j$  to be *adaptive*, so that they are chosen as part of the learning process. This leads to a consideration of multi-layer adaptive networks, as discussed in Chapters 4 and 5.

### 3.6 Fisher's linear discriminant

As the final topic of this chapter we consider a rather different approach to linear discriminants, introduced by Fisher (1936). In Section 1.4 we encountered the problem of the 'curse of dimensionality' whereby the design of a good classifier becomes rapidly more difficult as the dimensionality of the input space

increases. One way of dealing with this problem is to pre-process the data so as to reduce its dimensionality before applying a classification algorithm. The Fisher discriminant aims to achieve an optimal linear dimensionality reduction. It is therefore not strictly a discriminant itself, but it can easily be used to construct a discriminant. As well as being an important technique in its own right, the Fisher discriminant provides insight into the representations learned by multi-layer networks, as discussed in Section 6.6.1.

### 3.6.1 Two classes

One very simple approach to dimensionality reduction, motivated by our earlier discussion of single-layer networks, is to use a linear projection of the data onto a one-dimensional space, so that an input vector  $\mathbf{x}$  is projected onto a value  $y$  given by

$$y = \mathbf{w}^T \mathbf{x} \quad (3.77)$$

where, as before,  $\mathbf{w}$  is a vector of adjustable weight parameters. Note that this expression does not contain any bias parameter. We shall return to this point shortly. In general, the projection onto one dimension leads to a considerable loss of information, and classes which 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 which maximizes the class separation. To begin with, consider a two-class problem in which there are  $N_1$  points of class  $C_1$  and  $N_2$  points of class  $C_2$ . The mean vectors of the two classes are given by

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}^n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}^n. \quad (3.78)$$

We might think of defining the separation of the classes, when projected onto  $\mathbf{w}$ , as being the separation of the projected class means. This suggests that we might choose  $\mathbf{w}$  so as to maximize

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) \quad (3.79)$$

where

$$m_k = \mathbf{w}^T \mathbf{m}_k \quad (3.80)$$

is the class mean of the projected data from class  $C_k$ . 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 (Appendix C) 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,

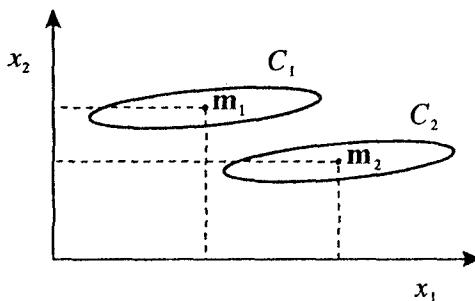


Figure 3.15. A schematic illustration of why it is important to take account of the within-class covariances when constructing the Fisher linear discriminant criterion. Projection of the data onto the  $x_1$ -axis leads to greater separation of the projected class means than does projection onto the  $x_2$ -axis, and yet it leads to greater class overlap. The problem is resolved by taking account of the within-class scatter of the data points.

however, as illustrated in Figure 3.15. This shows two classes which are well separated in the original two-dimensional space ( $x_1, x_2$ ). We see that projection onto the  $x_1$ -axis gives a much larger separation of the projected class means than does projection onto the  $x_2$ -axis. Nevertheless, separation of the projected data is much better when the data is projected onto the  $x_2$ -axis than when it is projected onto the  $x_1$ -axis. This difficulty arises from the substantial difference of the within-class spreads along the two axis directions. The resolution proposed by Fisher is to maximize a function which represents the difference between the projected class means, normalized by a measure of the within-class scatter along the direction of  $w$ .

The projection formula (3.77) transforms the set of labelled data points in  $\mathbf{x}$  into a labelled set in the one-dimensional space  $y$ . The within-class scatter of the transformed data from class  $C_k$  is described by the within-class covariance, given by

$$s_k^2 = \sum_{n \in C_k} (y^n - m_k)^2 \quad (3.81)$$

and we can define the total within-class covariance for the whole data set to be simply  $s_1^2 + s_2^2$ . We therefore arrive at the Fisher criterion given by

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}. \quad (3.82)$$

We can make the dependence on  $\mathbf{w}$  explicit by using (3.77), (3.80) and (3.81) to rewrite the Fisher criterion in the form

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (3.83)$$

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

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \quad (3.84)$$

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

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}^n - \mathbf{m}_1)(\mathbf{x}^n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}^n - \mathbf{m}_2)(\mathbf{x}^n - \mathbf{m}_2)^T. \quad (3.85)$$

Differentiating (3.83) with respect to  $\mathbf{w}$ , we find that  $J(\mathbf{w})$  is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}. \quad (3.86)$$

From (3.84) 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. Thus, we can drop any scalar factors. Multiplying both sides of (3.86) by  $\mathbf{S}_W^{-1}$  we then obtain

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1). \quad (3.87)$$

This 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. 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 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. In doing this we note that  $y = \mathbf{w}^T \mathbf{x}$  is the sum of a set of random variables, and so we may invoke the central limit theorem (see page 37) and model the class-conditional density functions  $p(y|\mathcal{C}_k)$  using normal distributions. The techniques of Chapter 2 can then be used to find the parameters of the normal distributions by maximum likelihood, and the formalism of Chapter 1 then gives an expression for the optimal threshold.

Once we have obtained a suitable weight vector and a threshold, the procedure for deciding the class of a new vector is identical to that of the perceptron network of Section 3.5. We can therefore view the Fisher criterion as a specific procedure for choosing the weights (and subsequently the bias) in a single-layer network. More conventionally, however, it is regarded as a technique for dimensionality reduction, a subject which is discussed at greater length in Chapter 8. In reducing the dimensionality of the data we are discarding information, and this cannot reduce (and will typically increase) the theoretical minimum achievable error rate. Dimensionality reduction may be worthwhile in practice, however, as

it alleviates problems associated with the curse of dimensionality. Thus, with finite-sized data sets, reduction of the dimensionality may well lead to overall improvements in the performance of a classifier system.

### 3.6.2 Relation to the least-squares approach

The least-squares approach to the determination of a linear discriminant was based on the goal of making the network outputs 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 taken the target values to be +1 for class  $C_1$  and -1 for class  $C_2$ . 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  $C_1$  to be  $N/N_1$ , where  $N_1$  is the number of patterns in class  $C_1$ , and  $N$  is the total number of patterns. This target value approximates the reciprocal of the prior probability for class  $C_1$ . For class  $C_2$  we shall take the targets to be  $-N/N_2$ .

The sum-of-squares error function can be written

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^n + w_0 - t^n)^2. \quad (3.88)$$

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

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^n + w_0 - t^n) = 0 \quad (3.89)$$

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^n + w_0 - t^n) \mathbf{x}^n = 0. \quad (3.90)$$

From (3.89), and making use of our choice of target coding scheme for the  $t^n$ , we obtain an expression for the bias in the form

$$w_0 = -\mathbf{w}^T \mathbf{m} \quad (3.91)$$

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

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^n = \frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2). \quad (3.92)$$

After some straightforward algebra, and again making use of the choice of  $t^n$ , the second equation (3.90) becomes

$$\left( \mathbf{S}_W + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_2 - \mathbf{m}_1) \quad (3.93)$$

where  $\mathbf{S}_W$  is defined by (3.85),  $\mathbf{S}_B$  is defined by (3.84), and we have substituted for the bias using (3.91). Using (3.84) we note that  $\mathbf{S}_B \mathbf{w}$  is always in the direction of  $(\mathbf{m}_2 - \mathbf{m}_1)$ . Thus we can write

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \quad (3.94)$$

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 (3.91). This tells us that a new vector  $\mathbf{x}$  should be classified as belonging to class  $C_1$  if  $\mathbf{w}^T(\mathbf{x} - \mathbf{m}) > 0$  and class  $C_2$  otherwise.

### 3.6.3 Several classes

We now consider the generalization of the Fisher discriminant to several classes, and we shall assume that the dimensionality of the input space is greater than the number of classes, so that  $d > c$ . Also, 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 rows of a matrix  $\mathbf{W}$ , so that

$$\mathbf{y} = \mathbf{W}\mathbf{x} \quad (3.95)$$

The generalization of the within-class covariance matrix to the case of  $c$  classes follows from (3.85) to give

$$\mathbf{S}_W = \sum_{k=1}^c \mathbf{S}_k \quad (3.96)$$

where

$$\mathbf{S}_k = \sum_{n \in C_k} (\mathbf{x}^n - \mathbf{m}_k)(\mathbf{x}^n - \mathbf{m}_k)^T \quad (3.97)$$

and

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}^n \quad (3.98)$$

where  $N_k$  is the number of patterns in class  $C_k$ . 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

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}^n - \mathbf{m})(\mathbf{x}^n - \mathbf{m})^T \quad (3.99)$$

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

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^n = \frac{1}{N} \sum_{k=1}^c N_k \mathbf{m}_k \quad (3.100)$$

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 (3.96) and (3.97), plus an additional matrix  $\mathbf{S}_B$  which we identify as a measure of the between-class covariance

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B \quad (3.101)$$

where

$$\mathbf{S}_B = \sum_{k=1}^c N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T. \quad (3.102)$$

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

$$\mathbf{s}_W = \sum_{k=1}^c \sum_{n \in C_k} (\mathbf{y}^n - \boldsymbol{\mu}_k)(\mathbf{y}^n - \boldsymbol{\mu}_k)^T \quad (3.103)$$

and

$$\mathbf{s}_B = \sum_{k=1}^c N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T \quad (3.104)$$

where

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}^n, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^c N_k \boldsymbol{\mu}_k. \quad (3.105)$$

Again we wish to construct a scalar which 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

$$J(\mathbf{W}) = \text{Tr} \{ \mathbf{s}_W^{-1} \mathbf{s}_B \} \quad (3.106)$$

where  $\text{Tr}\{\mathbf{M}\}$  denotes the trace of a matrix  $\mathbf{M}$ . This criterion can then be rewritten as an explicit function of the projection matrix  $\mathbf{W}$  in the form

$$J(\mathbf{W}) = \text{Tr} \{ (\mathbf{W} \mathbf{S}_W \mathbf{W}^T)^{-1} (\mathbf{W} \mathbf{S}_B \mathbf{W}^T) \}. \quad (3.107)$$

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  $\mathbf{S}_W^{-1} \mathbf{S}_B$  which correspond to the  $d'$  largest eigenvalues.

There is one important result which is common to all such criteria, which is worth emphasizing. We first note from (3.102) that  $\mathbf{S}_B$  is composed of the sum of  $c$  matrices, each of which is an outer product of two vectors and therefore of rank 1. In addition only  $(c - 1)$  of these matrices are independent as a result of the constraint (3.100). Thus,  $\mathbf{S}_B$  has rank at most equal to  $(c - 1)$  and so there are at most  $(c - 1)$  non-zero eigenvalues. This shows that the projection down onto the  $(c - 1)$ -dimensional sub-space 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  $(c - 1)$  linear 'features' by this means (Fukunaga, 1990). Dimensionality reduction and feature extraction are discussed at greater length in Chapter 8.

### Exercises

- 3.1 (\*)** Consider a point  $\hat{\mathbf{x}}$  which lies on the plane  $y(\hat{\mathbf{x}}) = 0$ , where  $y(\mathbf{x})$  is given by (3.1). By minimizing the distance  $\|\mathbf{x} - \hat{\mathbf{x}}\|$  with respect to  $\hat{\mathbf{x}}$  subject to this constraint, show that the value of the linear discriminant function  $y(\mathbf{x})$  gives a (signed) measure of the perpendicular distance  $L$  of the point  $\mathbf{x}$  to the decision boundary  $y(\mathbf{x}) = 0$  of the form

$$L = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}. \quad (3.108)$$

- 3.2 (\*)** There are several possible ways in which to generalize the concept of a linear discriminant function from two classes to  $c$  classes. One possibility would be to use  $(c - 1)$  linear discriminant functions, such that  $y_k(\mathbf{x}) > 0$  for inputs  $\mathbf{x}$  in class  $C_k$  and  $y_k(\mathbf{x}) < 0$  for inputs not in class  $C_k$ . By drawing a simple example in two dimensions for  $c = 3$ , show that this approach can lead to regions of  $\mathbf{x}$ -space for which the classification is ambiguous. Another approach would be to use one discriminant function  $y_{jk}(\mathbf{x})$  for each possible pair of classes  $C_j$  and  $C_k$ , such that  $y_{jk}(\mathbf{x}) > 0$  for patterns in class  $C_j$ , and  $y_{jk}(\mathbf{x}) < 0$  for patterns in class  $C_k$ . For  $c$  classes we would need  $c(c - 1)/2$  discriminant functions. Again, by drawing a specific example in two dimensions for  $c = 3$ , show that this approach can also lead to ambiguous regions.

- 3.3 (\*)** Consider a mixture model of the form (2.71) in which the component densities are given by

$$p(\mathbf{x}|j) = \prod_{i=1}^d P_{ji}^{x_i} (1 - P_{ji})^{1-x_i} \quad (3.109)$$

which is equivalent to (3.22). Show that the maximum likelihood solution for the parameters  $P_{ji}$  is given by

$$\hat{P}_{ji} = \frac{\sum_n P(j|\mathbf{x}^n)x_i^n}{\sum_n P(j|\mathbf{x}^n)} \quad (3.110)$$

where  $P(j|\mathbf{x})$  is the posterior probability for component  $j$  corresponding to an input vector  $\mathbf{x}$  and is given, from Bayes' theorem, by

$$P(j|\mathbf{x}) = \frac{p(\mathbf{x}|j)P(j)}{\sum_k p(\mathbf{x}|k)P(k)} \quad (3.111)$$

and  $P(j)$  is the corresponding prior probability.

- 3.4 (\*\*) 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**

$$\mathbf{x} = \sum_n \alpha_n \mathbf{x}^n \quad (3.112)$$

where  $\alpha_n \geq 0$  and  $\sum_n \alpha_n = 1$ . Consider a second set of points  $\{\mathbf{z}^n\}$  and its corresponding convex hull. 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  $\hat{\mathbf{w}}^T \mathbf{z}^n + w_0 < 0$  for all  $\mathbf{z}^n$ . 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.

- 3.5 (\*\*) Draw all  $2^2 = 4$  dichotomies of  $N = 2$  points in one dimension, and hence show that the fraction of such dichotomies which are linearly separable is 1.0. By considering the binomial expansion of  $2^d = (1+1)^d$ , verify that the summation in (3.30) does indeed give  $F = 1$  when  $N = d + 1$  for any  $d$ . Similarly, by drawing all  $2^4 = 16$  dichotomies of  $N = 4$  points in one dimension, show that the fraction of dichotomies which are linearly separable is 0.5. By considering the binomial expansion of  $2^{2d+1} = (1+1)^{2d+1}$ , show from (3.30) that the fraction of dichotomies which are linearly separable for  $N = 2(d + 1)$  is given by  $F(2d + 2, d) = 0.5$  for any  $N$ . Verify that these results are consistent with Figure 3.7.**

- 3.6 (\*\*\* Generate and plot a set of data points in two dimensions, drawn from two classes each of which is described by a Gaussian class-conditional density function. Implement the gradient descent algorithm for training a logistic discriminant, and plot the decision boundary at regular intervals**

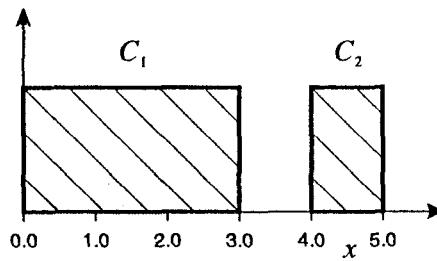


Figure 3.16. Distribution of data in one dimension drawn from two classes, used in Exercise 3.7.

during the training procedure on the same graph as the data. Explore the effects of choosing different values for the learning rate parameter  $\eta$ . Compare the behaviour of the sequential and batch weight update procedures described by (3.52) and (3.54) respectively.

- 3.7 (\*\*) Consider** data in one dimension drawn from two classes having the distributions shown in Figure 3.16. What is the ratio of the prior probabilities for the two classes? Find the linear discriminant function  $y(x) = wx + w_0$  which minimizes the sum-of-squares error function defined by

$$E = 3 \int_0^3 \{y(x) - 1\}^2 dx + \int_4^5 \{y(x) + 1\}^2 dx \quad (3.113)$$

where the target values are  $t = +1$  for class  $C_1$  and  $t = -1$  for class  $C_2$ . Show that the decision boundary given by  $y(x) = 0$  just fails to separate the two classes. Would a single-layer perceptron necessarily find a solution which separates the two classes exactly? Justify your answer. Discuss briefly the advantages and limitations of the least-squares and perceptron algorithms in the light of these results.

- 3.8 (\*) Prove** that, for arbitrary vectors  $w$  and  $\hat{w}$ , the following inequality is satisfied:

$$\frac{(\hat{w}^T w)^2}{\|\hat{w}\|^2 \|w\|^2} \leq 1. \quad (3.114)$$

Hence, using the results (3.73) and (3.76) from the proof of the perceptron convergence theorem given in the text, show that an upper limit on the number of weight updates needed for convergence of the perceptron algorithm is given by

$$\tau_{\max} = \frac{\|\hat{w}\|^2 \|\phi\|_{\max}^2}{\min_n (\hat{w}^T \phi^n t^n)^2}. \quad (3.115)$$

- 3.9 (\*\*\*)** Generate a data set consisting of a small number of vectors in two dimensions, each belonging to one of two classes. Write a numerical implementation of the perceptron learning algorithm, and plot both the data points and the decision boundary after every iteration. Explore the behaviour of the algorithm both for data sets which are linearly separable and for those which are not.
- 3.10 (\*)** Use a Lagrange multiplier (Appendix C) to show that, for two classes, the projection vector which maximizes the separation of the projected class means given by (3.79), subject to the constraint  $\|\mathbf{w}\|^2 = 1$ , is given by  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ .
- 3.11 (\*\*)** Using the definitions of the between-class and within-class covariance matrices given by (3.84) and (3.85) respectively, together with (3.91) and (3.92) and the choice of target values described in Section 3.6.2, show that the expression (3.90) which minimizes the sum-of-squares error function can be written in the form (3.93).
- 3.12 (\*)** By making use of (3.98), show that the total covariance matrix  $\mathbf{S}_T$  given by (3.99) can be decomposed into within-class and between-class covariance matrices as in (3.101), where the within-class covariance matrix  $\mathbf{S}_W$  is given by (3.96) and (3.97), and the between-class covariance matrix  $\mathbf{S}_B$  is given by (3.102).