

10

Linear Discrimination

In linear discrimination, we assume that instances of a class are linearly separable from instances of other classes. This is a discriminant-based approach that estimates the parameters of the linear discriminant directly from a given labeled sample.

10.1 Introduction

WE REMEMBER from the previous chapters that in classification we define a set of discriminant functions $g_j(\mathbf{x})$, $j = 1, \dots, K$, and then we

choose C_i if $g_i(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x})$

Previously, when we discussed methods for classification, we first estimated the prior probabilities, $\hat{P}(C_i)$, and the class likelihoods, $\hat{p}(\mathbf{x}|C_i)$, then used Bayes' rule to calculate the posterior densities. We then defined the discriminant functions in terms of the posterior, for example,

$$g_i(\mathbf{x}) = \log \hat{P}(C_i|\mathbf{x})$$

LIKELIHOOD-BASED
CLASSIFICATION

This is called *likelihood-based classification*, and we have previously discussed the parametric (chapter 5), semiparametric (chapter 7), and nonparametric (chapter 8) approaches to estimating the class likelihoods, $p(\mathbf{x}|C_i)$.

DISCRIMINANT-BASED
CLASSIFICATION

We are now going to discuss *discriminant-based classification* where we assume a model directly for the discriminant, bypassing the estimation of likelihoods or posteriors. The discriminant-based approach, as we also saw for the case of decision trees in chapter 9, makes an assumption on the form of the discriminant between the classes and makes no assumption about, or requires no knowledge of the densities—for example,

whether they are Gaussian, or whether the inputs are correlated, and so forth.

We define a model for the discriminant

$$g_i(\mathbf{x}|\Phi_i)$$

explicitly parameterized with the set of parameters Φ_i , as opposed to a likelihood-based scheme that has implicit parameters in defining the likelihood densities. This is a different inductive bias: instead of making an assumption on the form of the class densities, we make an assumption on the form of the boundaries separating classes.

Learning is the optimization of the model parameters Φ_i to maximize the quality of the separation, that is, the classification accuracy on a given labeled training set. This differs from the likelihood-based methods that search for the parameters that maximize sample likelihoods, separately for each class.

In the discriminant-based approach, we do not care about correctly estimating the densities inside class regions; all we care about is the correct estimation of the *boundaries* between the class regions. Those who advocate the discriminant-based approach (e.g., Vapnik 1995) state that estimating the class densities is a harder problem than estimating the class discriminants, and it does not make sense to solve a hard problem to solve an easier problem. This is of course true only when the discriminant can be approximated by a simple function.

In this chapter, we concern ourselves with the simplest case where the discriminant functions are linear in \mathbf{x} :

$$(10.1) \quad g_i(\mathbf{x}|\mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} = \sum_{j=1}^d w_{ij}x_j + w_{i0}$$

LINEAR DISCRIMINANT

The *linear discriminant* is used frequently mainly due to its simplicity: both the space and time complexities are $\mathcal{O}(d)$. The linear model is easy to understand: the final output is a weighted sum of the input attributes x_j . The magnitude of the weight w_j shows the importance of x_j and its sign indicates if the effect is positive or negative. Most functions are additive in that the output is the sum of the effects of several attributes where the weights may be positive (enforcing) or negative (inhibiting). For example, when a customer applies for credit, financial institutions calculate the applicant's credit score that is generally written as a sum of the effects of various attributes; for example, yearly income has a positive effect (higher incomes increase the score).

In many applications, the linear discriminant is also quite accurate. We know, for example, that when classes are Gaussian with a shared covariance matrix, the optimal discriminant is linear. The linear discriminant, however, can be used even when this assumption does not hold, and the model parameters can be calculated without making any assumptions on the class densities. We should always use the linear discriminant before trying a more complicated model to make sure that the additional complexity is justified.

As always, we formulate the problem of finding a linear discriminant function as a search for the parameter values that minimize an error function. In particular, we concentrate on *gradient* methods for optimizing a criterion function.

10.2 Generalizing the Linear Model

QUADRATIC
DISCRIMINANT

When a linear model is not flexible enough, we can use the *quadratic discriminant* function and increase complexity

$$(10.2) \quad g_i(\mathbf{x}|\mathbf{W}_i, \mathbf{w}_i, w_{i0}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

but this approach is $O(d^2)$ and we again have the bias/variance dilemma: the quadratic model, though is more general, requires much larger training sets and may overfit on small samples.

HIGHER-ORDER TERMS
PRODUCT TERMS

An equivalent way is to preprocess the input by adding *higher-order terms*, also called *product terms*. For example, with two inputs x_1 and x_2 , we can define new variables

$$z_1 = x_1, z_2 = x_2, z_3 = x_1^2, z_4 = x_2^2, z_5 = x_1 x_2$$

and take $\mathbf{z} = [z_1, z_2, z_3, z_4, z_5]^T$ as the input. The linear function defined in the five-dimensional \mathbf{z} space corresponds to a nonlinear function in the two-dimensional \mathbf{x} space. Instead of defining a nonlinear function (discriminant or regression) in the original space, what we do is to define a suitable nonlinear transformation to a new space where the function can be written in a linear form.

We write the discriminant as

$$(10.3) \quad g_i(\mathbf{x}) = \sum_{j=1}^k w_j \phi_{ij}(\mathbf{x})$$

BASIS FUNCTION

where $\phi_{ij}(\mathbf{x})$ are *basis functions*. Higher-order terms are only one set of possible basis functions; other examples are

- $\sin(x_1)$
- $\exp(-(x_1 - m)^2/c)$
- $\exp(-\|\mathbf{x} - \mathbf{m}\|^2/c)$
- $\log(x_2)$
- $1(x_1 > c)$
- $1(ax_1 + bx_2 > c)$

POTENTIAL FUNCTION

where m, a, b, c are scalars, \mathbf{m} is a d -dimensional vector, and $1(b)$ returns 1 if b is true and returns 0 otherwise. The idea of writing a nonlinear function as a linear sum of nonlinear basis functions is an old idea and was originally called *potential functions* (Aizerman, Braverman, and Rozonoer 1964). Multilayer perceptrons (chapter 11) and radial basis functions (chapter 12) have the advantage that the parameters of the basis functions can be fine-tuned to the data during learning. In chapter 13, we discuss support vector machines that use kernel functions built from such basis functions.

10.3 Geometry of the Linear Discriminant

10.3.1 Two Classes

Let us start with the simpler case of two classes. In such a case, one discriminant function is sufficient:

$$\begin{aligned}
 g(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\
 &= (\mathbf{w}_1^T \mathbf{x} + w_{10}) - (\mathbf{w}_2^T \mathbf{x} + w_{20}) \\
 &= (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (w_{10} - w_{20}) \\
 &= \mathbf{w}^T \mathbf{x} + w_0
 \end{aligned}$$

and we

$$\text{choose } \begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

WEIGHT VECTOR
THRESHOLD

This defines a hyperplane where \mathbf{w} is the *weight vector* and w_0 is the *threshold*. This latter name comes from the fact that the decision rule can be rewritten as follows: choose C_1 if $\mathbf{w}^T \mathbf{x} > -w_0$, and choose C_2

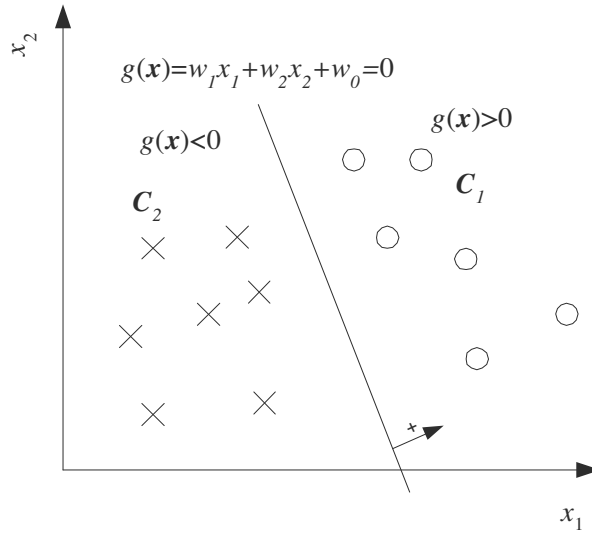


Figure 10.1 In the two-dimensional case, the linear discriminant is a line that separates the examples from two classes.

otherwise. The hyperplane divides the input space into two half-spaces: the decision region \mathcal{R}_1 for C_1 and \mathcal{R}_2 for C_2 . Any \mathbf{x} in \mathcal{R}_1 is on the *positive* side of the hyperplane and any \mathbf{x} in \mathcal{R}_2 is on its *negative* side. When \mathbf{x} is $\mathbf{0}$, $g(\mathbf{x}) = w_0$ and we see that if $w_0 > 0$, the origin is on the positive side of the hyperplane, and if $w_0 < 0$, the origin is on the negative side, and if $w_0 = 0$, the hyperplane passes through the origin (see figure 10.1).

Take two points \mathbf{x}_1 and \mathbf{x}_2 both on the decision surface; that is, $g(\mathbf{x}_1) = g(\mathbf{x}_2) = 0$, then

$$\begin{aligned}\mathbf{w}^T \mathbf{x}_1 + w_0 &= \mathbf{w}^T \mathbf{x}_2 + w_0 \\ \mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) &= 0\end{aligned}$$

and we see that \mathbf{w} is normal to any vector lying on the hyperplane. Let us rewrite \mathbf{x} as (Duda, Hart, and Stork 2001)

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

where \mathbf{x}_p is the normal projection of \mathbf{x} onto the hyperplane and r gives us the distance from \mathbf{x} to the hyperplane, negative if \mathbf{x} is on the negative

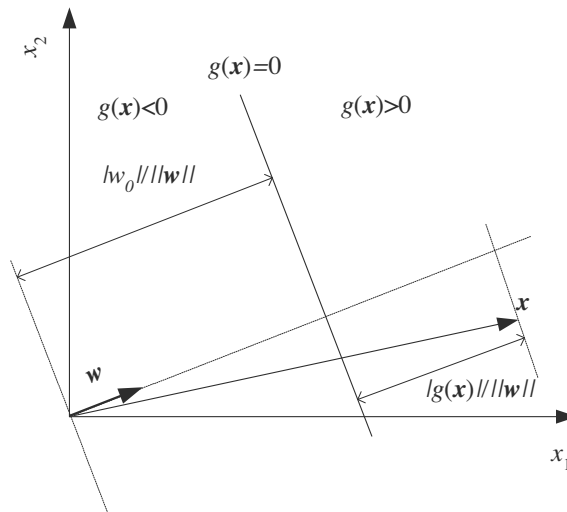


Figure 10.2 The geometric interpretation of the linear discriminant.

side, and positive if \mathbf{x} is on the positive side (see figure 10.2). Calculating $g(\mathbf{x})$ and noting that $g(\mathbf{x}_p) = 0$, we have

$$(10.4) \quad r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

We see then that the distance to origin is

$$(10.5) \quad r_0 = \frac{w_0}{\|\mathbf{w}\|}$$

Thus w_0 determines the location of the hyperplane with respect to the origin, and \mathbf{w} determines its orientation.

10.3.2 Multiple Classes

When there are $K > 2$ classes, there are K discriminant functions. When they are linear, we have

$$(10.6) \quad g_i(\mathbf{x}|\mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

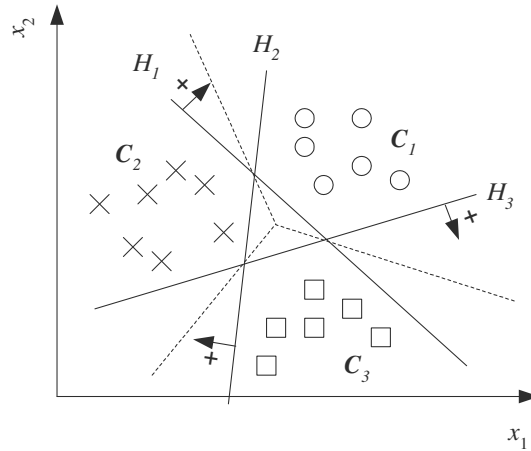


Figure 10.3 In linear classification, each hyperplane H_i separates the examples of C_i from the examples of all other classes. Thus for it to work, the classes should be linearly separable. Dotted lines are the induced boundaries of the linear classifier.

We are going to talk about learning later on but for now, we assume that the parameters, \mathbf{w}_i, w_{i0} , are computed so as to have

$$(10.7) \quad g_i(\mathbf{x}|\mathbf{w}_i, w_{i0}) = \begin{cases} > 0 & \text{if } \mathbf{x} \in C_i \\ \leq 0 & \text{otherwise} \end{cases}$$

LINEARLY SEPARABLE
CLASSES

for all \mathbf{x} in the training set. Using such discriminant functions corresponds to assuming that all classes are *linearly separable*; that is, for each class C_i , there exists a hyperplane H_i such that all $\mathbf{x} \in C_i$ lie on its positive side and all $\mathbf{x} \in C_j, j \neq i$ lie on its negative side (see figure 10.3).

During testing, given \mathbf{x} , ideally, we should have only one $g_j(\mathbf{x}), j = 1, \dots, K$ greater than 0 and all others should be less than 0, but this is not always the case: The positive half-spaces of the hyperplanes may overlap, or, we may have a case where all $g_j(\mathbf{x}) < 0$. These may be taken as *reject* cases, but the usual approach is to assign \mathbf{x} to the class having the highest discriminant:

$$(10.8) \quad \text{Choose } C_i \text{ if } g_i(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x})$$

Remembering that $|g_i(\mathbf{x})|/\|\mathbf{w}_i\|$ is the distance from the input point to the hyperplane, assuming that all \mathbf{w}_i have similar length, this assigns the

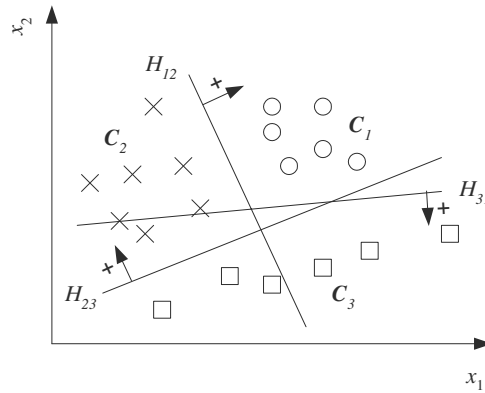


Figure 10.4 In pairwise linear separation, there is a separate hyperplane for each pair of classes. For an input to be assigned to C_1 , it should be on the positive side of H_{12} and H_{13} (which is the negative side of H_{31}); we do not care for the value of H_{23} . In this case, C_1 is not linearly separable from other classes but is pairwise linearly separable.

LINEAR CLASSIFIER

point to the class (among all $g_j(\mathbf{x}) > 0$) to whose hyperplane the point is most distant. This is called a *linear classifier*, and geometrically it divides the feature space into K convex decision regions \mathcal{R}_i (see figure 10.3).

10.4 Pairwise Separation

PAIRWISE SEPARATION

If the classes are not linearly separable, one approach is to divide it into a set of linear problems. One possibility is *pairwise separation* of classes (Duda, Hart, and Stork 2001). It uses $K(K-1)/2$ linear discriminants, $g_{ij}(\mathbf{x})$, one for every pair of distinct classes (see figure 10.4):

$$g_{ij}(\mathbf{x} | \mathbf{w}_{ij}, w_{ij0}) = \mathbf{w}_{ij}^T \mathbf{x} + w_{ij0}$$

The parameters $\mathbf{w}_{ij}, j \neq i$ are computed during training so as to have

$$(10.9) \quad g_{ij}(\mathbf{x}) = \begin{cases} > 0 & \text{if } \mathbf{x} \in C_i \\ \leq 0 & \text{if } \mathbf{x} \in C_j \\ \text{don't care} & \text{otherwise} \end{cases} \quad i, j = 1, \dots, K \text{ and } i \neq j$$

that is, if $\mathbf{x}^t \in C_k$ where $k \neq i, k \neq j$, then \mathbf{x}^t is not used during training of $g_{ij}(\mathbf{x})$.

During testing, we

choose C_i if $\forall j \neq i, g_{ij}(\mathbf{x}) > 0$

In many cases, this may not be true for any i and if we do not want to reject such cases, we can relax the conjunction by using a summation and choosing the maximum of

$$(10.10) \quad g_i(\mathbf{x}) = \sum_{j \neq i} g_{ij}(\mathbf{x})$$

Even if the classes are not linearly separable, if the classes are pairwise linearly separable—which is much more likely—pairwise separation can be used, leading to nonlinear separation of classes (see figure 10.4). This is another example of breaking down a complex (e.g., nonlinear) problem, into a set of simpler (e.g., linear) problems. We have already seen decision trees (chapter 9) that use this idea, and we will see more examples of this in chapter 17 on combining multiple models, for example, error-correcting output codes, and mixture of experts, where the number of linear models is less than $\mathcal{O}(K^2)$.

10.5 Parametric Discrimination Revisited

In chapter 5, we saw that if the class densities, $p(\mathbf{x}|C_i)$, are Gaussian and share a common covariance matrix, the discriminant function is linear

$$(10.11) \quad g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

where the parameters can be analytically calculated as

$$(10.12) \quad \begin{aligned} \mathbf{w}_i &= \Sigma^{-1} \boldsymbol{\mu}_i \\ w_{i0} &= -\frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \log P(C_i) \end{aligned}$$

Given a dataset, we first calculate the estimates for $\boldsymbol{\mu}_i$ and Σ and then plug the estimates, \mathbf{m}_i , \mathbf{S} , in equation 10.12 and calculate the parameters of the linear discriminant.

Let us again see the special case where there are two classes. We define $y \equiv P(C_1|\mathbf{x})$ and $P(C_2|\mathbf{x}) = 1 - y$. Then in classification, we

$$\text{choose } C_1 \text{ if } \begin{cases} y > 0.5 \\ \frac{y}{1-y} > 1 \\ \log \frac{y}{1-y} > 0 \end{cases} \quad \text{and } C_2 \text{ otherwise}$$

LOGIT
LOG ODDS $\log y/(1 - y)$ is known as the *logit* transformation or *log odds* of y . In the case of two normal classes sharing a common covariance matrix, the log odds is linear:

$$\begin{aligned}
 \text{logit}(P(C_1|\mathbf{x})) &= \log \frac{P(C_1|\mathbf{x})}{1 - P(C_1|\mathbf{x})} = \log \frac{P(C_1|\mathbf{x})}{P(C_2|\mathbf{x})} \\
 &= \log \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} + \log \frac{P(C_1)}{P(C_2)} \\
 &= \log \frac{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)]}{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)]} + \log \frac{P(C_1)}{P(C_2)} \\
 (10.13) \quad &= \mathbf{w}^T \mathbf{x} + w_0
 \end{aligned}$$

where

$$\begin{aligned}
 \mathbf{w} &= \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\
 (10.14) \quad w_0 &= -\frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)^T \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \log \frac{P(C_1)}{P(C_2)}
 \end{aligned}$$

The inverse of logit

$$\log \frac{P(C_1|\mathbf{x})}{1 - P(C_1|\mathbf{x})} = \mathbf{w}^T \mathbf{x} + w_0$$

LOGISTIC
SIGMOID is the *logistic* function, also called the *sigmoid* function (see figure 10.5):

$$(10.15) \quad P(C_1|\mathbf{x}) = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]}$$

During training, we estimate $\mathbf{m}_1, \mathbf{m}_2, \mathbf{S}$ and plug these estimates in equation 10.14 to calculate the discriminant parameters. During testing, given \mathbf{x} , we can either

1. calculate $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(\mathbf{x}) > 0$, or
2. calculate $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0)$ and choose C_1 if $y > 0.5$,

because $\text{sigmoid}(0) = 0.5$. In this latter case, sigmoid transforms the discriminant value to a posterior probability. This is valid when there are two classes and one discriminant; we see in section 10.7 how we can estimate posterior probabilities for $K > 2$.

10.6 Gradient Descent

In likelihood-based classification, the parameters were the sufficient statistics of $p(\mathbf{x}|C_i)$ and $P(C_i)$, and the method we used to estimate the parameters is maximum likelihood. In the discriminant-based approach,

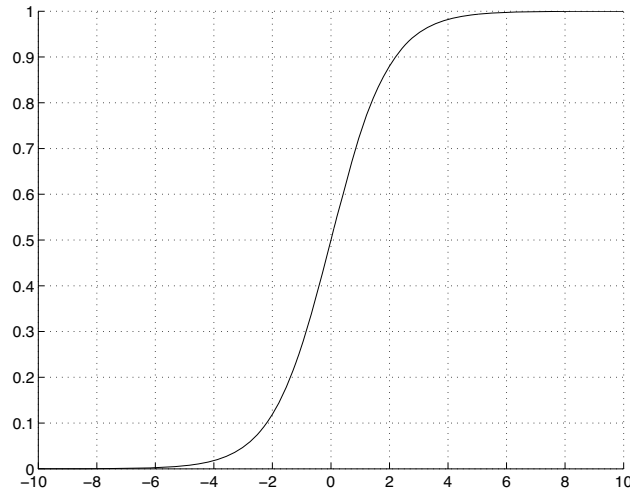


Figure 10.5 The logistic, or sigmoid, function.

the parameters are those of the discriminants, and they are optimized to minimize the classification error on the training set. When \mathbf{w} denotes the set of parameters and $E(\mathbf{w}|X)$ is the error with parameters \mathbf{w} on the given training set X , we look for

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w}|X)$$

In many cases, some of which we will see shortly, there is no analytical solution and we need to resort to iterative optimization methods, the most commonly employed being that of *gradient descent*. When $E(\mathbf{w})$ is a differentiable function of a vector of variables, we have the *gradient vector* composed of the partial derivatives

$$\nabla_{\mathbf{w}} E = \left[\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_d} \right]^T$$

and the *gradient descent* procedure to minimize E starts from a random \mathbf{w} , and at each step, updates \mathbf{w} , in the opposite direction of the gradient

$$(10.16) \quad \Delta w_i = -\eta \frac{\partial E}{\partial w_i}, \forall i$$

$$(10.17) \quad w_i = w_i + \Delta w_i$$

where η is called the *stepsize*, or *learning factor*, and determines how much to move in that direction. Gradient ascent is used to maximize a

GRADIENT DESCENT
GRADIENT VECTOR

function and goes in the direction of the gradient. When we get to a minimum (or maximum), the derivative is 0 and the procedure terminates. This indicates that the procedure finds the nearest minimum that can be a local minimum, and there is no guarantee of finding the global minimum unless the function has only one minimum. The use of a good value for η is also critical; if it is too small, the convergence may be too slow, and a large value may cause oscillations and even divergence.

Throughout this book, we use gradient methods that are simple and quite effective. We keep in mind, however, that once a suitable model and an error function is defined, the optimization of the model parameters to minimize the error function can be done by using one of many possible techniques. There are second-order methods and conjugate gradient that converge faster, at the expense of more memory and computation. More costly methods like simulated annealing and genetic algorithms allow a more thorough search of the parameter space and do not depend as much on the initial point.

10.7 Logistic Discrimination

10.7.1 Two Classes

LOGISTIC
DISCRIMINATION

In *logistic discrimination*, we do not model the class-conditional densities, $p(\mathbf{x}|C_i)$, but rather their ratio. Let us again start with two classes and assume that the log likelihood ratio is linear:

$$(10.18) \quad \log \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} = \mathbf{w}^T \mathbf{x} + w_0^o$$

This indeed holds when the class-conditional densities are normal (equation 10.13). But logistic discrimination has a wider scope of applicability; for example, \mathbf{x} may be composed of discrete attributes or may be a mixture of continuous and discrete attributes.

Using Bayes' rule, we have

$$\begin{aligned} \text{logit}(P(C_1|\mathbf{x})) &= \log \frac{P(C_1|\mathbf{x})}{1 - P(C_1|\mathbf{x})} \\ &= \log \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} + \log \frac{P(C_1)}{P(C_2)} \\ (10.19) \quad &= \mathbf{w}^T \mathbf{x} + w_0 \end{aligned}$$

where

$$(10.20) \quad w_0 = w_0^o + \log \frac{P(C_1)}{P(C_2)}$$

Rearranging terms, we get the sigmoid function again:

$$(10.21) \quad y = \hat{P}(C_1|\mathbf{x}) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]}$$

as our estimator of $P(C_1|\mathbf{x})$.

Let us see how we can learn \mathbf{w} and w_0 . We are given a sample of two classes, $\mathcal{X} = \{\mathbf{x}^t, r^t\}$, where $r^t = 1$ if $\mathbf{x} \in C_1$ and $r^t = 0$ if $\mathbf{x} \in C_2$. We assume r^t , given \mathbf{x}^t , is Bernoulli with probability $y^t \equiv P(C_1|\mathbf{x}^t)$ as calculated in equation 10.21:

$$r^t|\mathbf{x}^t \sim \text{Bernoulli}(y^t)$$

Here, we see the difference from the likelihood-based methods where we modeled $p(\mathbf{x}|C_i)$; in the discriminant-based approach, we model directly $r|\mathbf{x}$. The sample likelihood is

$$(10.22) \quad l(\mathbf{w}, w_0|\mathcal{X}) = \prod_t (y^t)^{(r^t)} (1 - y^t)^{(1-r^t)}$$

We know that when we have a likelihood function to maximize, we can always turn it into an error function to be minimized as $E = -\log l$, and in our case, we have *cross-entropy*:

CROSS-ENTROPY

$$(10.23) \quad E(\mathbf{w}, w_0|\mathcal{X}) = - \sum_t r^t \log y^t + (1 - r^t) \log(1 - y^t)$$

Because of the nonlinearity of the sigmoid function, we cannot solve directly and we use gradient descent to minimize cross-entropy, equivalent to maximizing the likelihood or the log likelihood. If $y = \text{sigmoid}(a) = 1/(1 + \exp(-a))$, its derivative is given as

$$\frac{dy}{da} = y(1 - y)$$

and we get the following update equations:

$$(10.24) \quad \begin{aligned} \Delta w_j &= -\eta \frac{\partial E}{\partial w_j} = \eta \sum_t \left(\frac{r^t}{y^t} - \frac{1 - r^t}{1 - y^t} \right) y^t (1 - y^t) x_j^t \\ &= \eta \sum_t (r^t - y^t) x_j^t, j = 1, \dots, d \\ \Delta w_0 &= -\eta \frac{\partial E}{\partial w_0} = \eta \sum_t (r^t - y^t) \end{aligned}$$

```

For  $j = 0, \dots, d$ 
   $w_j \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
  For  $j = 0, \dots, d$ 
     $\Delta w_j \leftarrow 0$ 
  For  $t = 1, \dots, N$ 
     $o \leftarrow 0$ 
    For  $j = 0, \dots, d$ 
       $o \leftarrow o + w_j x_j^t$ 
     $y \leftarrow \text{sigmoid}(o)$ 
    For  $j = 0, \dots, d$ 
       $\Delta w_j \leftarrow \Delta w_j + (r^t - y)x_j^t$ 
  For  $j = 0, \dots, d$ 
     $w_j \leftarrow w_j + \eta \Delta w_j$ 
Until convergence

```

Figure 10.6 Logistic discrimination algorithm implementing gradient descent for the single output case with two classes. For w_0 , we assume that there is an extra input x_0 , which is always $+1$: $x_0^t \equiv +1, \forall t$.

It is best to initialize w_j with random values close to 0; generally they are drawn uniformly from the interval $[-0.01, 0.01]$. The reason for this is that if the initial w_j are large in magnitude, the weighted sum may also be large and may saturate the sigmoid. We see from figure 10.5 that if the initial weights are close to 0, the sum will stay in the middle region where the derivative is nonzero and an update can take place. If the weighted sum is large in magnitude (smaller than -5 or larger than $+5$), the derivative of the sigmoid will be almost 0 and weights will not be updated.

Pseudocode is given in figure 10.6. We see an example in figure 10.7 where the input is one-dimensional. Both the line $w x + w_0$ and its value after the sigmoid are shown as a function of learning iterations. We see that to get outputs of 0 and 1, the sigmoid hardens, which is achieved by increasing the magnitude of w , or $\|w\|$ in the multivariate case.

Once training is complete and we have the final w and w_0 , during testing, given x^t , we calculate $y^t = \text{sigmoid}(w^T x^t + w_0)$ and we choose C_1 if $y^t > 0.5$ and choose C_2 otherwise. This implies that to minimize the number of misclassifications, we do not need to continue learning un-

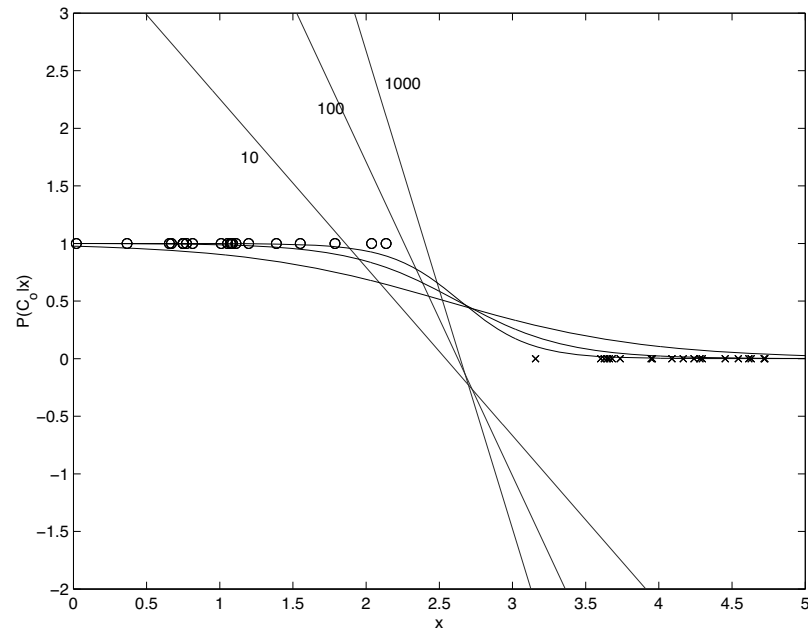


Figure 10.7 For a univariate two-class problem (shown with ‘ \circ ’ and ‘ \times ’), the evolution of the line $wx + w_0$ and the sigmoid output after 10, 100, and 1,000 iterations over the sample.

til all y^t are 0 or 1, but only until y^t are less than or greater than 0.5, that is, on the correct side of the decision boundary. If we do continue training beyond this point, cross-entropy will continue decreasing ($|w_j|$ will continue increasing to harden the sigmoid), but the number of misclassifications will not decrease. Generally, we continue training until the number of misclassifications does not decrease (which will be 0 if the classes are linearly separable). Actually *stopping early* before we have 0 training error is a form of regularization. Because we start with weights almost 0 and they move away as training continues, stopping early corresponds to a model with more weights close to 0 and effectively fewer parameters.

EARLY STOPPING

Note that though we assumed the log ratio of the class densities are linear to derive the discriminant, we estimate directly the posterior and never explicitly estimate $p(\mathbf{x}|C_i)$ or $P(C_i)$.

10.7.2 Multiple Classes

Let us now generalize to $K > 2$ classes. We take one of the classes, for example, C_K , as the reference class and assume that

$$(10.25) \quad \log \frac{p(\mathbf{x}|C_i)}{p(\mathbf{x}|C_K)} = \mathbf{w}_i^T \mathbf{x} + w_{i0}^o$$

Then we have

$$(10.26) \quad \frac{P(C_i|\mathbf{x})}{P(C_K|\mathbf{x})} = \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]$$

with $w_{i0} = w_{i0}^o + \log P(C_i)/P(C_K)$.

We see that

$$(10.27) \quad \begin{aligned} \sum_{i=1}^{K-1} \frac{P(C_i|\mathbf{x})}{P(C_K|\mathbf{x})} &= \frac{1 - P(C_K|\mathbf{x})}{P(C_K|\mathbf{x})} = \sum_{i=1}^{K-1} \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}] \\ &\Rightarrow P(C_K|\mathbf{x}) = \frac{1}{1 + \sum_{i=1}^{K-1} \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]} \end{aligned}$$

and also that

$$(10.28) \quad \begin{aligned} \frac{P(C_i|\mathbf{x})}{P(C_K|\mathbf{x})} &= \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}] \\ &\Rightarrow P(C_i|\mathbf{x}) = \frac{\exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]}{1 + \sum_{j=1}^{K-1} \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}]}, \quad i = 1, \dots, K-1 \end{aligned}$$

To treat all classes uniformly, we can write

$$(10.29) \quad y_i = \hat{P}(C_i|\mathbf{x}) = \frac{\exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]}{\sum_{j=1}^K \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}]}, \quad i = 1, \dots, K$$

SOFTMAX which is called the *softmax* function (Bridle 1990). If the weighted sum for one class is sufficiently larger than for the others, after it is boosted through exponentiation and normalization, its corresponding y_i will be close to 1 and the others will be close to 0. Thus it works like taking a maximum, except that it is differentiable; hence the name softmax. Softmax also guarantees that $\sum_i y_i = 1$.

Let us see how we can learn the parameters. In this case of $K > 2$ classes, each sample point is a multinomial trial with one draw; that is, $\mathbf{r}^t | \mathbf{x}^t \sim \text{Mult}_K(1, \mathbf{y}^t)$, where $y_i^t \equiv P(C_i | \mathbf{x}^t)$. The sample likelihood is

$$(10.30) \quad l(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = \prod_t \prod_i (y_i^t)^{r_i^t}$$

and the error function is again cross-entropy:

$$(10.31) \quad E(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = - \sum_t \sum_i r_i^t \log y_i^t$$

We again use gradient descent. If $y_i = \exp(a_i) / \sum_j \exp(a_j)$, we have

$$(10.32) \quad \frac{\partial y_i}{\partial a_j} = y_i(\delta_{ij} - y_j)$$

where δ_{ij} is the Kronecker delta, which is 1 if $i = j$ and 0 if $i \neq j$ (exercise 3). Given that $\sum_i r_i^t = 1$, we have the following update equations, for $j = 1, \dots, K$

$$\begin{aligned} \Delta \mathbf{w}_j &= \eta \sum_t \sum_i \frac{r_i^t}{y_i^t} y_i^t (\delta_{ij} - y_j^t) \mathbf{x}^t \\ &= \eta \sum_t \sum_i r_i^t (\delta_{ij} - y_j^t) \mathbf{x}^t \\ &= \eta \sum_t \left[\sum_i r_i^t \delta_{ij} - y_j^t \sum_i r_i^t \right] \mathbf{x}^t \\ &= \eta \sum_t (r_j^t - y_j^t) \mathbf{x}^t \\ (10.33) \quad \Delta w_{j0} &= \eta \sum_t (r_j^t - y_j^t) \end{aligned}$$

Note that because of the normalization in softmax, \mathbf{w}_j and w_{j0} are affected not only by $\mathbf{x}^t \in C_j$ but also by $\mathbf{x}^t \in C_i, i \neq j$. The discriminants are updated so that the correct class has the highest weighted sum after softmax, and the other classes have their weighted sums as low as possible. Pseudocode is given in figure 10.8. For a two-dimensional example with three classes, the contour plot is given in figure 10.9, and the discriminants and the posterior probabilities in figure 10.10.

During testing, we calculate all $y_k, k = 1, \dots, K$ and choose C_i if $y_i = \max_k y_k$. Again we do not need to continue training to minimize cross-entropy as much as possible; we train only until the correct class has the highest weighted sum, and therefore we can stop training earlier by checking the number of misclassifications.

When data are normally distributed, the logistic discriminant has a comparable error rate to the parametric, normal-based linear discriminant (McLachlan 1992). Logistic discrimination can still be used when the class-conditional densities are nonnormal or when they are not unimodal, as long as classes are linearly separable.

```

For  $i = 1, \dots, K$ , For  $j = 0, \dots, d$ ,  $w_{ij} \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
  For  $i = 1, \dots, K$ , For  $j = 0, \dots, d$ ,  $\Delta w_{ij} \leftarrow 0$ 
  For  $t = 1, \dots, N$ 
    For  $i = 1, \dots, K$ 
       $o_i \leftarrow 0$ 
      For  $j = 0, \dots, d$ 
         $o_i \leftarrow o_i + w_{ij} x_j^t$ 
      For  $i = 1, \dots, K$ 
         $y_i \leftarrow \exp(o_i) / \sum_k \exp(o_k)$ 
      For  $i = 1, \dots, K$ 
        For  $j = 0, \dots, d$ 
           $\Delta w_{ij} \leftarrow \Delta w_{ij} + (r_i^t - y_i) x_j^t$ 
      For  $i = 1, \dots, K$ 
        For  $j = 0, \dots, d$ 
           $w_{ij} \leftarrow w_{ij} + \eta \Delta w_{ij}$ 
  Until convergence

```

Figure 10.8 Logistic discrimination algorithm implementing gradient descent for the case with $K > 2$ classes. For generality, we take $x_0 \equiv 1, \forall t$.

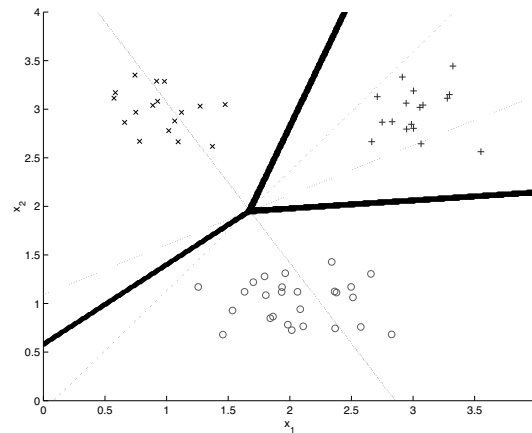


Figure 10.9 For a two-dimensional problem with three classes, the solution found by logistic discrimination. Thin lines are where $g_i(\mathbf{x}) = 0$, and the thick line is the boundary induced by the linear classifier choosing the maximum.

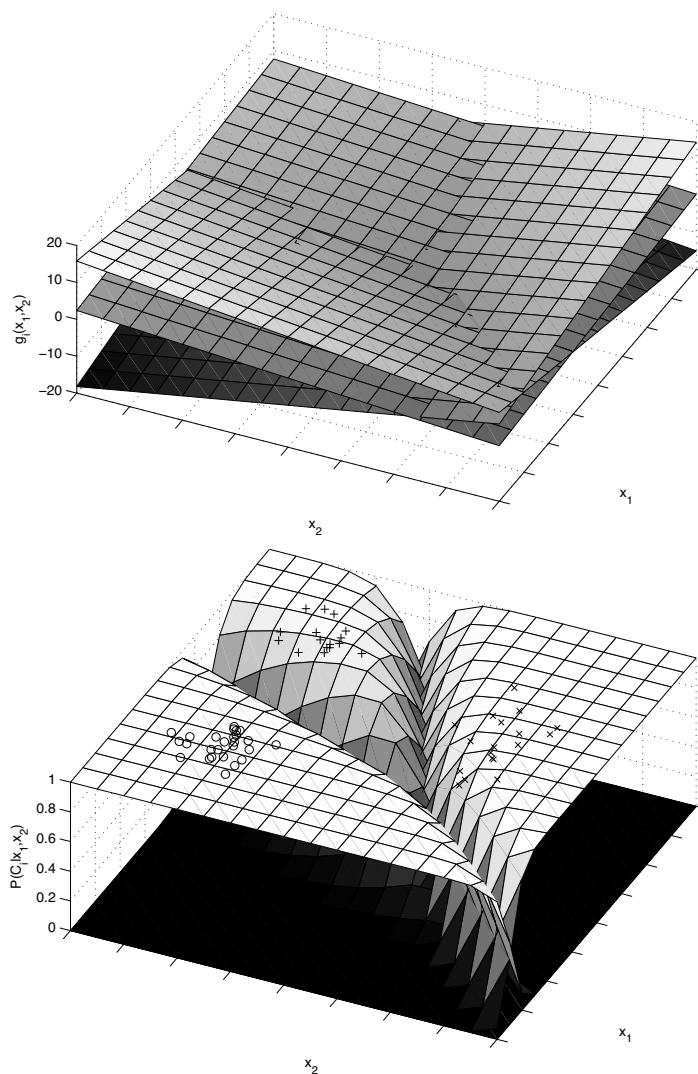


Figure 10.10 For the same example in figure 10.9, the linear discriminants (top), and the posterior probabilities after the softmax (bottom).

The ratio of class-conditional densities is of course not restricted to be linear (Anderson 1982; McLachlan 1992). Assuming a quadratic discriminant, we have

$$(10.34) \quad \log \frac{p(\mathbf{x}|C_i)}{p(\mathbf{x}|C_K)} = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

corresponding to and generalizing parametric discrimination with multivariate normal class-conditionals having different covariance matrices. When d is large, just as we can simplify (regularize) Σ_i , we can equally do it on \mathbf{W}_i by taking only its leading eigenvectors into account.

As discussed in section 10.2, any specified function of the basic variables can be included as \mathbf{x} -variates. One can, for example, write the discriminant as a linear sum of nonlinear basis functions

$$(10.35) \quad \log \frac{p(\mathbf{x}|C_i)}{p(\mathbf{x}|C_K)} = \mathbf{w}_i^T \boldsymbol{\phi}(\mathbf{x}) + w_{i0}$$

where $\boldsymbol{\phi}(\cdot)$ are the basis functions, which can be viewed as transformed variables. In neural network terminology, this is called a *multilayer perceptron* (chapter 11), and sigmoid is the most popular basis function. When a Gaussian basis function is used, the model is called *radial basis functions* (chapter 12). We can even use a completely nonparametric approach, for example, Parzen windows (chapter 8).

10.8 Discrimination by Regression

In regression, the probabilistic model is

$$(10.36) \quad r^t = y^t + \epsilon$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$. If $r^t \in \{0, 1\}$, y^t can be constrained to lie in this range using the sigmoid function. Assuming a linear model and two classes, we have

$$(10.37) \quad y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t + w_0) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x}^t + w_0)]}$$

Then the sample likelihood in regression, assuming $r|\mathbf{x} \sim \mathcal{N}(y, \sigma^2)$, is

$$(10.38) \quad l(\mathbf{w}, w_0|\mathcal{X}) = \prod_t \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(r^t - y^t)^2}{2\sigma^2} \right]$$

Maximizing the log likelihood is minimizing the sum of square errors:

$$(10.39) \quad E(\mathbf{w}, w_0 | \mathcal{X}) = \frac{1}{2} \sum_t (r^t - y^t)^2$$

Using gradient descent, we get

$$(10.40) \quad \begin{aligned} \Delta \mathbf{w} &= \eta \sum_t (r^t - y^t) y^t (1 - y^t) \mathbf{x}^t \\ \Delta w_0 &= \eta \sum_t (r^t - y^t) y^t (1 - y^t) \end{aligned}$$

This method can also be used when there are $K > 2$ classes. The probabilistic model is

$$(10.41) \quad \mathbf{r}^t = \mathbf{y}^t + \boldsymbol{\epsilon}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}_K(0, \sigma^2 \mathbf{I}_K)$. Assuming a linear model for each class, we have

$$(10.42) \quad y_i^t = \text{sigmoid}(\mathbf{w}_i^T \mathbf{x}^t + w_{i0}) = \frac{1}{1 + \exp[-(\mathbf{w}_i^T \mathbf{x}^t + w_{i0})]}$$

Then the sample likelihood is

$$(10.43) \quad l(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = \prod_t \frac{1}{(2\pi)^{K/2} |\Sigma|^{1/2}} \exp \left[-\frac{\|\mathbf{r}^t - \mathbf{y}^t\|^2}{2\sigma^2} \right]$$

and the error function is

$$(10.44) \quad E(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = \frac{1}{2} \sum_t \|\mathbf{r}^t - \mathbf{y}^t\|^2 = \frac{1}{2} \sum_t \sum_i (r_i^t - y_i^t)^2$$

The update equations for $i = 1, \dots, K$, are

$$(10.45) \quad \begin{aligned} \Delta \mathbf{w}_i &= \eta \sum_t (r_i^t - y_i^t) y_i^t (1 - y_i^t) \mathbf{x}^t \\ \Delta w_{i0} &= \eta \sum_t (r_i^t - y_i^t) y_i^t (1 - y_i^t) \end{aligned}$$

But note that in doing so, we do not make use of the information that only one of y_i needs to be 1 and all others are 0, or that $\sum_i y_i = 1$. The softmax function of equation 10.29 allows us to incorporate this extra information we have due to the outputs' estimating class posterior probabilities. Using sigmoid outputs in $K > 2$ case, we treat y_i as if they are independent functions.

Note also that for a given class, if we use the regression approach, there will be updates until the right output is 1 and all others are 0. This is not

in fact necessary because during testing, we are just going to choose the maximum anyway; it is enough to train only until the right output is larger than others, which is exactly what the softmax function does.

So this approach with multiple sigmoid outputs is more appropriate when the classes are *not* mutually exclusive and exhaustive. That is, for an \mathbf{x}^t , all r_i^t may be 0; namely, \mathbf{x}^t does not belong to any of the classes, or more than one r_i^t may be 1, when classes overlap.

10.9 Notes

The linear discriminant, due to its simplicity, is the classifier most used in pattern recognition (Duda, Hart, and Stork 2001; McLachlan 1992). We discussed the case of Gaussian distributions with a common covariance matrix in chapter 4 and Fisher's linear discriminant in chapter 6, and in this chapter we discuss the logistic discriminant. In chapter 11, we discuss the perceptron that is the neural network implementation of the linear discriminant. In chapter 13, we will discuss support vector machines, another type of linear discriminant.

Logistic discrimination is covered in more detail in Anderson 1982 and in McLachlan 1992. Logistic (sigmoid) is the inverse of logit, which is the *canonical link* in case of Bernoulli samples. Softmax is its generalization to multinomial samples. More information on such *generalized linear models* is given in McCulloch and Nelder 1989.

Generalizing linear models by using nonlinear basis functions is a very old idea. We will discuss multilayer perceptrons (chapter 11) and radial basis functions (chapter 12) where the parameters of the basis functions can also be learned from data while learning the discriminant. Support vector machines (chapter 13) use kernel functions built from such basis functions.

GENERALIZED LINEAR
MODELS

10.10 Exercises

1. For each of the following basis functions, describe where it is nonzero:
 - a. $\sin(x_1)$
 - b. $\exp(-(x_1 - a)^2/c)$
 - c. $\exp(-\|\mathbf{x} - \mathbf{a}\|^2/c)$
 - d. $\log(x_2)$

- e. $1(x_1 > c)$
- f. $1(ax_1 + bx_2 > c)$
- 2. For the two-dimensional case of figure 10.2, show equations 10.4 and 10.5.
- 3. Show that the derivative of the softmax, $y_i = \exp(a_i) / \sum_j \exp(a_j)$, is $\partial y_i / \partial a_j = y_i(\delta_{ij} - y_j)$, where δ_{ij} is 1 if $i = j$ and 0 otherwise.
- 4. With $K = 2$, show that using two softmax outputs is equal to using one sigmoid output.
- 5. How can we learn \mathbf{W}_i in equation 10.34?
- 6. In using quadratic (or higher-order) discriminants as in equation 10.34, how can we keep variance under control?
- 7. What is the implication of the use of a single η for all x_j in gradient descent?
- 8. In the univariate case for classification as in figure 10.7, what do w and w_0 correspond to?
- 9. Let us say for univariate x , $x \in (2, 4)$ belong to C_1 and $x < 2$ or $x > 4$ belong to C_2 . How can we separate the two classes using a linear discriminant?

10.11 References

- Aizerman, M. A., E. M. Braverman, and L. I. Rozonoer. 1964. "Theoretical Foundations of the Potential Function Method in Pattern Recognition Learning." *Automation and Remote Control* 25: 821-837.
- Anderson, J. A. 1982. "Logistic Discrimination." In *Handbook of Statistics*, Vol. 2, *Classification, Pattern Recognition and Reduction of Dimensionality*, ed. P. R. Krishnaiah and L. N. Kanal, 169-191. Amsterdam: North Holland.
- Bridle, J. S. 1990. "Probabilistic Interpretation of Feedforward Classification Network Outputs with Relationships to Statistical Pattern Recognition." In *Neurocomputing: Algorithms, Architectures and Applications*, ed. F. Fogelman-Soulie and J. Herault, 227-236. Berlin: Springer.
- Duda, R. O., P. E. Hart, and D. G. Stork. 2001. *Pattern Classification*, 2nd ed. New York: Wiley.
- McCullagh, P., and J. A. Nelder. 1989. *Generalized Linear Models*. London: Chapman and Hall.
- McLachlan, G. J. 1992. *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- Vapnik, V. 1995. *The Nature of Statistical Learning Theory*. New York: Springer.