4 Linear Models Chapter for Classification XER Input vector Goal : assign of K discrete classes CK, K=1...K × to one decision space is divided into regions (Pt) the input So are called decision boundaries decision surfaces boundaries whose or for classification mean that the decision boundaries Linear models vector functions र्भ input linear are Ж 2.5 $2.0 \cdot$ 1.5(i.e. D-1 dimensional hyperplane) 1.0 v_2

Chapter 4 Linear Models for Classification

-0.5

-1.0 + -1.0

-0.5

0.0

0.5

 x_1

1.0

1.5

2.0

Probablistic models

Two class problem (binary representation) single target variable $t \in \{0, 1\}$ s.t t=1 represents class C_1 and t=0 represents class C_2 The value of t can be interpreted as the probability that class is C_1

K>2 classes problem (multiclass) K-dimensional vector #: one hot vector (1-of-K coding) If the class is C;, then all elements the of the are zero except t_i , $t_i = 1$. (1, 0, 0) (0, 1, 0)Again we can interprete the as the probability that the class is Ck.

In the linear model, the model prediction y(x,w) was by a linear function of W. given For classification problem, we need to predict discrete class more generally posterior probabilities. or => generalize the model in which we transform the linear function of wr using a nonlinear function $f(\cdot)$ so that $Y(x) = f(w^{T}x + w_{0})$

f(·) is known as an activation function.⁴ Its inverse is called a link function⁴ Chapter 4 Linear Models for Classification⁴

The decision boundaries correspond to YCX) = constant (i.e. with the = constant) and hence decision boundaries are linear function of x. In contrast to the models used for regression, classifications not linear in wr due to f(.) are As regression models, we can use a fixed nonlinear transformation with a vector of basis functions $\mathbb{Q}(\mathbb{X})$. We begin by considering classification directly in the original input space X.

4.1 Discriminant Functions Discriminant function takes X and assigns it one of K classes Ck We restrict attention to linear discriminants (decision boundaries are hyperplanes) 4.1.1 Two classes

The simplest linear discriminant function

$$YC$$
 := W^T + W_o

where W_{i} is the weight vector and W_{0} is a bias. An input K is assigned to C_{1} if $Y(K) \ge 0$ and is assigned to C_{2} if Y(K) < 0 \Rightarrow decision boundary is defined by Y(K) = 0. $M_{i}^{T} \times + W_{0}$

Arbitrary point x and let x_{\perp} be its orthogonal projection onto decision surface so that $x = x_{\perp} + r \frac{w_{\perp}}{\|w_{\parallel}\|}$

where

 $r = \frac{\gamma c^{(K)}}{\|W\|}$

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



We can use dummy input $x_0 = 1$ and then define $\widehat{W}_1 := (W_0, W_1)$ and $\widehat{X}_1 := (X_0, X_1)$ so that $Y(X_1) = \widehat{W}_1^T \widehat{X}_1$

In this case the decision boundaries are D-dimensional hyperplanes passing through the origin of D+1 dim input space.

4.1.2 Multiple classes

Now consider the extension of linear discriminants to K>2 classes

One-versus-the-rest classifier Use K-1 classifiers each of which solves a two-class problem separating class CK from other class This method leads to regions that are ambiguously classified



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

Consider a single K-class discriminant comprising. K linear
functions of the form

$$Y_k(x) := W_k^T x + W_{k0}$$
 $W_k \in \mathbb{R}^D$ weight vector
 $W_k o \in \mathbb{R}$ bias
Assign a point x to class C_k if $Y_k(x) > Y_j(x)$
 $\forall j \neq k$.
So the decision boundary between C_k and C_j is given
by
 $(W_k - W_j)^T x + (W_{k0} - W_{j0}) = 0$.

i.e. D-1 dimensional hyperplane.

Chapter 4 Linear Models for Classification







Since the discriminant functions are linear, we obtain

$$Y_{k}(x) = \pi Y_{k}(x_{A}) + (1 - \pi) Y_{k}(x_{B})$$

Because
$$X_A$$
 and X_B lie inside R_E , $Y_E(X_A) > Y_j(X_A)$ and
 $Y_E(X_B) > Y_j(X_B)$ $\forall j \neq K$, hence $Y_E(\hat{X}) > Y_j(\hat{X})$ $\forall j \neq K$

4.1.3 Least squares for classification Consider a classification problem with K classes with 1-of-k scheme for the target vector #. The minimization of SSE function is the method that it approximates the conditional expectation E[#1*] of the target values given the input X.

Each class C_k , k=1,2,..., K $Y_k C_k = W_k^T + W_{k0} = (\widetilde{W_k}^T)$

We can group these linear models using vector notation

$$Y(x) := \widetilde{W}^T \widetilde{X} = \begin{pmatrix} \chi(x) \\ \vdots \\ \chi_k(x) \end{pmatrix} \quad k-dim$$

 $M_{k}^{2} \widetilde{X}$

where \tilde{X} is the augmented input vector $(1, \tilde{X})^T$ and $(DH) \times K$ $\tilde{(DH)} \times K$

$$\widetilde{W}_{l} = \left(\widetilde{W}_{l}, \widetilde{W}_{2}, ..., \widetilde{W}_{k} \right) \qquad \widetilde{W}_{k} = \left(W_{k0}, W_{k} \right)^{\prime}$$

Input x is assigned to the class for which the output

$$Y_E = \widetilde{W}_E^T \mathscr{X} (= W_E^T \mathscr{X} + W_{EO})$$
 is largest
Determine the parameter matrix \widetilde{W} by minimizing a SSE.
Consider a training tata set $\{\mathscr{X}n, tln\}$ $n=1, ...N$
 $\mathscr{X}n \in \mathbb{R}^D$, $tln \in \mathbb{R}^k$ (one hot vector)
 $\mathscr{X}n \in \mathbb{R}^D$, $tln \in \mathbb{R}^k$ (one hot vector)
Define a matrix T whose n^{th} row is the vector tln
 $\mathscr{X} x(DT)$
 \mathscr{X} whose n^{th} row is the vector $\widetilde{\mathscr{X}}_n^T = (1...\mathfrak{X}_n)^T$
 $\widetilde{SSE} = \frac{1}{2} \sum_{n=1}^{N} || \widetilde{W}_n^T \widetilde{\mathscr{X}}_n - tln ||^2$
Chapter Linear Models for Classification 17

SSE function can be written as

$$E_{P}(\widehat{W}) = \frac{1}{2} \operatorname{Tr} \left\{ (\widehat{X} \widehat{W} - \mathbb{T})^{\mathsf{T}} (\widehat{X} \widehat{W} - \mathbb{T}) \right\}$$

Set the gradient w.r.t \widetilde{W} to zero vector. So we obtain the minimizing solution of $E_{\mathcal{D}}(\widetilde{W})$ for \widetilde{W} as solution $\widetilde{W} = (\widetilde{\chi}^T \widetilde{\chi})^{-1} \widetilde{\chi}^T T = \widetilde{\chi}^+ T$

where X^{\dagger} is the pseudo-inverse of \overline{X} . The discriminant function is given by $Y(X) = \overline{W}^{\intercal} \overline{X} = \overline{T}^{\intercal} (\overline{X}^{\dagger})^{\intercal} \overline{X}$.

W: the parameter matrix whose kth column is
$$W_{k}$$

X: the matrix whose n^{th} row is K_{n}^{T}
Then, $E_{D}(\overline{W}) = \frac{1}{2} \operatorname{Tr} \int (X W + 1 W_{0}^{T} - T)^{T} (X W + 1 W_{0}^{T} - T)$
where $k - \dim 1 := (1, 1, ..., 1)^{T}$, $W_{0} := (W_{10}, W_{20}, ..., W_{k0})^{T}$
Calculate the derivative of $E_{D}(W)$ w.r.t W_{0}
 $\nabla_{W_{0}} E_{D}(\overline{W}) = 2N W_{0} + 2(X W - T)^{T} 1$
 $K_{X1} M_{XK} K_{X1}$
(See later!)
Charter 4 linear Models for Classification

We have obtained the discriminant function using least square approach. $\gamma(\mathcal{K}) = \widetilde{\mathcal{W}}^T \widetilde{\mathcal{K}} = \pi^T (\widetilde{\mathcal{K}}^+)^T \widetilde{\mathcal{K}}$

where $\widetilde{\mathcal{X}}^{\dagger} = (\widetilde{\mathcal{X}}^{\mathsf{T}} \widetilde{\mathcal{X}})^{-1} \widetilde{\mathcal{X}}^{\mathsf{T}}.$

This discriminant function does not have any probabilistic interpretation and is not robust to outliers (least square)



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



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

Chapter 4 Linear Models for Classification

Least square approach

regression

4.1.4 Fisher's Linear discriminant Consider a linear classification in terms of dimensionality reduction Input XGIR. Consider a projection to one timension using $y := W^T$ Threshold we on y. So if $y \ge w_0$ then x is classified as class C1 otherwise class C2 Considerable loss of information and overlapping in one dimension

Goal: determine w/ or select projection maximizing the class separation.

Consider a two classes problem with N_1 points of class C_1 and N_2 points of class of C_2

The mean vectors of the two classes

$$ml_{i} := \frac{1}{N_{i}} \sum_{n \in C_{i}} \#_{n} \qquad ml_{2} := \frac{1}{N_{2}} \sum_{n \in C_{2}} \#_{n}$$

First choose w to maximize the difference of projected means $M_2 - M_1 = W_1^T (M_1 - M_1)$ where $M_k := W_1^T M_k$ Chapter 4 Linear Models for Classification 23

We constrain
$$w$$
 to have unit length, i.e. $\sum_{k} w_{k}^{2} = 1$
Using a Lagrange multiplier, we find
 $w/ oC (m_{1} - m_{1})$ (see Fig. 4.6)
Second, consider a small variance within each class
The within - class variance of the transformed (projected
data from class C_{k} is given by
 $S_{k}^{2} := \sum_{n \in C_{k}} (Y_{n} - m_{k})^{2}$

where $y_n = W^T \times n$ Chapter 4 Linear Models for Classification

24

The Fisher criterion: maximize

$$(4,26) \qquad J(w):=\frac{(m_2-m_1)^2}{s_1^2+s_2^2} \qquad between-class variance total within-class variance total within-cl$$

$$J(w) = \frac{\|w^{T}(m_{12} - m_{11})\|^{2}}{\sum_{n \in C_{1}} (w^{T} x_{n} - m_{11})^{2} + \sum_{n \in C_{2}} (w^{T} x_{n} - m_{11})^{2}}$$

$$\|w^{T}(m_{1} - m_{1})\|^{2} = [w^{T}(m_{2} - m_{1})][w^{T}(m_{2} - m_{1})]^{T} = w^{T}S_{B}w$$

where $S_{B} := (m_{12} - m_{11})(m_{12} - m_{11})^{T}$ (4.27) $DX_{1} = [XD]$

$$S_{1}^{2} + S_{2}^{2} = \sum_{n \in C_{1}} \left[W^{T} (\mathscr{K}_{n} - M_{1}) \right]^{2} + \sum_{n \in C_{2}} \left[W^{T} (\mathscr{K}_{n} - M_{1}) \right]^{2}$$
$$= W^{T} S_{w}^{1} W' + W^{T} S_{w}^{2} W' = W^{T} S_{w} W'$$

where
$$S_w^{\kappa} := \sum_{n \in G_{\kappa}} (M_n - M_{\kappa}) (M_n - M_{\kappa})^T$$
, $K = 1$, 2

$$S_{w} := \sum_{c_{1}} (\mathscr{K}_{n} - \mathsf{M}_{i}) (\mathscr{K}_{n} - \mathsf{M}_{i})^{T} + \sum_{c_{2}} (\mathscr{K}_{n} - \mathsf{M}_{2}) (\mathscr{K}_{n} - \mathsf{M}_{2})^{T} \quad (4.28)$$

Thus
$$J(w) = \frac{w_T S_B w_T}{w_T S_w w_T}$$
 S_B : between - class covariance matrix
 S_w : within - class covariance matrix
 W D-dim IXP DXD DXI

Differentiating
$$J(W) = W'' =$$

Multiplying both sides of (4.29) by $5w^{-1}$, we then obtain $w/ C S_w^{-1} (m_1 - m_1)$ (4.30)

Note that if within - class covariance is isotropic $(S_w = \pi I)$ then solution W is proportional to $M_2 - M_1$ (4.30) is known as Fisher's linear discriminant. This is the direction for projection of the data down to I-dimension.



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

4.1.5 Relation to least squares
Two approaches of linear discriminants for two-class problem
The least squares make the model predictions as close
as possible as to a set of target values.
Minimize
$$\sum_{n=1}^{N} || \frac{w_n^T \hat{x}_n}{y_{L} \hat{x}_n} - t_n ||^2$$

The Fisher criterion was derived by requiring maximum class separation in the 1-dim output space

Chapter 4 Linear Models for Classification



The sum-of-squares error function can be written

$$E = \frac{1}{2} \sum_{n=1}^{N} (w^{T} x_{n} + w_{o} - t_{n})^{2}$$

$$\mathscr{K}_{n} \in \mathbb{R}^{D}, \quad w = (w_{1}, \dots, w_{p})^{T}$$

Set the derivative of E w.r.t wo and W/ to zero,

$$\frac{\partial E}{\partial w_0} = \sum_{n=1}^{N} (w_n^T X_n + w_0 - t_n) = 0$$

$$\nabla_{w_{1}} E = \sum_{n=1}^{N} (w_{1}^{T} *_{n} + w_{o} - t_{n}) *_{n} = 0 \qquad (4.33)$$

Thus,
$$w_0 = -w_1^T m_1$$
 (4.34)

where m is the mean of total input data set and is given by

$$\mathbf{m}_{1} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{X}_{n} = \frac{1}{N} (N_{1} \mathbf{m}_{1} + N_{2} \mathbf{m}_{2})$$

To obtain (4.34) we have used

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

Chapter 4 Linear Models for Classification

.

(4.33) can be written (Excercise 4.6) $(S_{11} + \frac{N_1 N_2}{N_1} S_B) W = N CM_1 - Ml_2)$ is defined by C4.28) and SB is defined by where Sw (4.20)Since SBW is always in direction of (m2-m1), we can write $W_{1} \circ C S_{W_{1}}^{-1} (m_{12} - m_{1})$

where we have ignored irrelevant scale factors.



4.1.6 Fisher's discriminant for multiple classes
$$(k > 2)$$

WLOCT, assume $D > k$
The generalization of within - covariance matrix to the case
of k classes follows from $(4:28)$ to give
 $S_w = \sum_{k=1}^{k} S_k$ (input space)

where

$$S_{k} := \sum_{n \in C_{k}} (x_{n} - m_{k}) (x_{n} - m_{k})^{T}$$

$$M_{k} = \frac{1}{N_{k}} \sum_{n \in C_{k}} \mathscr{K}_{n}$$
where N_F is is of patterns in C_F

Consider the total covariance matrix

$$S_{T} := \sum_{n=1}^{N} (X_{n} - m) (X_{n} - m)^{T} \quad (input space)$$

where m_1 is the mean of the total tota set. This S_T can be tecomposed into the sum of the within - class covariance matrix S_W and an additional matrix S_B

 $S_{+} = S_{+} + S_{R}$ **Chapter 4 Linear Models for Classification**

We identify SB as a measure of the between-class covariance

(4.46)
$$S_{B} = \sum_{k=1}^{k} N_{k} (m_{k} - m) (m_{k} - m)^{T} (input space)$$

where
$$k = 1, .., D'$$
. (D-dim weight w_k)

The weight vectors 1 Wk1 can be considered to be the

columns of a matrix W/ (DXD') so that

$$W^{T} = W^{T} \times D' - \dim$$

Chapter 4 Linear Models for Classification

Now define similar matrices in the project $D'-\dim y - space$ $S_{w} = \sum_{k=1}^{K} \sum_{n \in C_{k}} (Y_{n} - M_{k}) (Y_{n} - M_{k})^{T} \quad (D' \times D' \text{ matrix})$

and

$$S_{B} = \sum_{k=1}^{K} N_{k} (M_{k} - M) (M_{k} - M)^{T} \quad (D' \times D' \text{ matrix})$$

where

To determine
$$W'$$
, we need to define a scalar (benefit)
which is large when S_B is large and when S_W is small.
Consider $J(W) := Tr (S_W^{-1} S_B)$ footure
space

This criterion can be written as an explicit function in the form $J(W) = Tr f (W^T S_W W)^{-1} (W^T S_B W) f$ $D \times P D \times P D \times P'$ in put space Remark

(4.46) (def of S_B), S_B is the sum of K matrices - From and each of which is of rank 1. - Because of the definition of m, only (K-1) of these matrices are independent - Thus SB has rank at most (k-1) and so there are at most (K-1) eigenvalues. - So the projection onto the (K-1) dim subspace spanned by the eigenvectors of SB does not change J(W). More than $(k - 1)^{\text{Chapter 4 Linear Models for Classification}}$ are meaningless 41

4.1.7 The Perceptron algorithm (linear discriminant model) Two-class classification

X input vector

 $\overline{\Phi}(\mathbf{x})$ its feature vector for a fixed nonlinear function $\overline{\Phi}(\cdot)$. Linear model of the form

 $\gamma(x) := f(w^T \overline{\phi}(x))$ parameter vector W

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

Chapter 4 Linear Models for Classification

$$f(\alpha) := \begin{cases} 1 & \alpha \ge 0 \\ -1 & \alpha < 0 & 0 \end{cases}$$

Here $\Phi(x)$ include a bias component $\beta_0(x) = 1$. For the perceptron, target values t=1 for C_1 , t=-1 for C_2 How to determine w? to define error function of W?? How Misclassification rate?

Perceptron criterion
Idea: if
$$\#n$$
 is in class C_1 , then $\#^T \overline{\Phi}(\#n) > 0$
 $\#$ $\#$ C_2 , then $\#^T \overline{\Phi}(\#n) \ll 0$
 $t_n = -1$
Using $t \in \{-1, 1\}$ target coding, we are seeking $\#$ set
 $\#^T \overline{\Phi}(\#)$ to > 0

The perceptron criterion is given by

$$E_{P}(w) := -\sum w^{T} \overline{\mathcal{Q}}(w) t_{n} \qquad (4.54)$$

where *M* denotes the set of all misclassified patterns.
So the total error function is piecewise linear for *w*.
If *x* is correctly classified then the contribution to the
error is zero.
The stochastic gradient descent algorithm to this error
$$W_{i}^{(C+1)}$$
 := $W_{i}^{(C)} - \eta \nabla E_{p}(w) = W_{i}^{(C)} + \eta \bar{\Phi}(w_{n}) t_{n}$
where η is the learning rate parameter. (put $\eta = 1$)

If pattern is correctly classified then the W/ remains unchanges In 计 is incorrectly classified, case $\overline{\Phi}(X_n)$ add estimate onto the for current D W/ 更CKn) while subtract from for C₂, Wr. se







Figure 4.7 Illustration of the convergence of the perceptron learning algorithm, showing data points from two classes (red and blue) in a two-dimensional feature space (ϕ_1, ϕ_2) . The top left plot shows the initial parameter vector w shown as a black arrow together with the corresponding decision boundary (black line), in which the arrow points towards the decision region which classified as belonging to the red class. The data point circled in green is misclassified and so its feature vector is added to the current weight vector, giving the new decision boundary shown in the top right plot. The bottom left plot shows the next misclassified point to be considered, indicated by the green circle, and its feature vector is again added to the weight vector giving the decision boundary shown in the bottom right plot for which all data points are correctly classified.

Remark

- In view of (4.54) and (4.55), the contribution to the error from a misclassified pattern will be reduced $-w' \quad \Phi(x_n) t_n = -w' \quad \Phi(x_n) t_n - (\Phi(x_n) t_n)^T \Phi(x_n) t_n$ single component of Ep $< - W^{(C)T} \overline{\Phi}(x_n) t_n$ where we have set n = 1 and used $|| \Phi(x_n) t_n || > 0$ - This does not the contribution to the error function from all misclassified patterns (other)

- The change in Wr may have caused some previously correctly classified patterns to become misclassified that the training data set is linearly separable, - In case perceptron learning algorithm is guarrenteed to fint an exact solution in a finite number of steps (by perceptron convergence theorem) - Perceptron does not provide probabilistic output Can generalize K>2 classes not on linear combinations of based basis fixed functions **Chapter 4 Linear Models for Classification**

48



Probablistic models

- generative
- discriminative

$$p(x(c_i)) \rightarrow 2d$$

 $p(x)$
 $p(x)$

4.2 Probabilistic Generative model
Discriminative and generative approaches to classification.
Consider the case of two classes.

$$P(C_1 | \mathscr{K}) = \frac{P(\mathscr{K} | C_1) P(C_1)}{P(\mathscr{K} | C_1) P(C_1) + P(\mathscr{K} | C_2) P(C_2)}$$

 $= \frac{1}{1 + exp(-\alpha)} = e^{-\alpha}$
(4.57)
 (4.57)
where we have defined

$$a = l_n \frac{P(X(C_i) P(C_i)}{(4.58)}$$

Chapter 4 Rineal Models for Classification

Remark of sigmoid

- Bounded function
- Symmetry property $\sigma(-a) = 1 \sigma(a)$
- The inverse of the logistic sigmoid is given by

$$a = ln\left(\frac{\sigma}{1-\sigma}\right)$$
 logit function

K > 2 classes

$$PCC_{E}(x) = \frac{P(x|C_{E}) PCC_{E}}{\Sigma_{j} P(x|C_{j}) PCC_{j}} = \frac{\exp(Ca_{E})}{\Sigma_{j} \exp(a_{i})}$$
softmax function
which is known as the normalized exponential (multiclass
generalization of the logistic sigmoid. Here a_{E} are defined
by
$$a_{E} := ln(P(x|C_{E}) PCC_{E})).$$

4.2.1 Continuous inputs X e IR^D continuous vector $P(X | C_{E})$ the class - conditional densities are Gaussian and Assume all classes share the same covariance matrix. (only different) I.e. $P(X|C_{k}) = \frac{1}{(2\pi)^{N_{1}}} \frac{1}{|\Sigma|^{N_{2}}} \exp\left(-\frac{1}{2}(X - M_{k})^{T} \Sigma^{-1}(X - M_{k})\right)$

Here Σ is independent of class C_{k} .

Consider the case of two classes. From
$$(4.5^{n})$$
 and (4.5^{n}) ,

$$P(C_{1}|k) = \sigma(w_{1}^{T} + w_{0})$$

$$= \sigma(a)$$

$$A = ln \frac{P(k|C_{1}) P(C_{1})}{P(k|C_{2}) P(C_{2})}$$
where we have defined
$$Wi := \sum^{-1} (M_{1} - M_{2})$$

$$w_{0} := -\frac{1}{2} M_{1}^{T} \sum^{-1} M_{1} + \frac{1}{2} M_{2}^{T} \sum^{-1} M_{2} + ln \frac{P(C_{1})}{P(C_{2})}$$
Because of the assumption of common covariance matrices, it
becomes a linear function of x_{1} in the argument of the
logistic sigmoid.
Chapter 4 Linear Models for Classification

Thus, decision boundary $(\times st p CC_{k} | \star) = C)$ is a linear function of \star . The prior $p CC_{k}$ enter only through the bias parameter

Wo.

For the general case of
$$K > 2$$
 classes under the assumption
shared covariance matrix of $p C \times I C_{k}$)
 $a_{k} (x) := W_{k}^{T} \times + W_{k0}$
where we have defined
 $a_{k} = ln(p C \times I C_{k}) p C_{k})$
 $W_{k} := \Sigma^{T} (M_{k})$
 $W_{k} := \Sigma^{T} (M_{k})$
 $w_{k0} := -\frac{1}{2} M_{k}^{T} \Sigma^{T} M_{k} + ln P C C_{k})$
 $a_{k} (x)$ is again linear function of x .

Chapter 4 Linear Models for Classification

If	each	class -	conditional	density	PCXICE)	has	its ou	'n
Covar	iance	matrix	Σ_{k} , the	n the	concellations	ත්	quadrati	C
form	0	F ×	will no la	onger occ	cur.			
So	we	obtain	quadratic	functions	र्ठ 🗶	givin	g rise	
to	a	quaratic	dis cri mi nant					



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

4.2.2 Maximum likelihood solution Two classes classification Dota set 1%, tn1, n=1...N, $\%_n \in \mathbb{R}^D$, tn = 1 or 0 tn = 1 denotes class C_1 and t=0 denotes class C_2 Graussian class conditional density with a shared covariance matrix

Denote the prior class probability
$$p(C_1) = \pi$$
, so that $p(C_2)$
= $1 - \pi$.
**n from class C_1 , $t_n = 1$ hence
 $P(*n, C_1) = p(C_1) p(*n | C_1) = \pi N((*n | M_1, \Sigma))$
Similarly, for class C_2 , $t_n = 0$
 $P(*n, C_2) = p(C_2) p(*n | C_2) = (1 - \pi) N((*n | M_2, \Sigma))$

Thus the likelihood function is given by

$$P(\#, X \mid \pi, M_1, M_2, \Sigma) = \prod_{n=1}^{N} [\pi N (X_n \mid M_1, \Sigma)]^{tn} [(1-\pi) N (X_n \mid M_2, \Sigma)]^{1-tn}$$

where
$$\# = (t_1, ..., t_N)^T$$
, $\chi = (\chi_1, \chi_2, ..., \chi_N)^T$

As usual, we maximize the log of the likelihood function. Consider first π . The log likelihood function of π is

$$\sum_{n=1}^{N} \int t_n \ln \pi + (1-t_n) \ln (1-\pi) f$$

Setting the derivative writ π equal to 0. So we obtain

$$\pi = \frac{1}{N} \sum_{n=1}^{N} t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} \qquad C_1 \forall l_2$$

where N_1 (resp. N_2) is # of points in C_1 (resp. C_2) Thus MLE for π is simply the fraction of points in C_1 Now consider the maximization w.r.t M_1 . The terms of loglikelihood function depending on M_1 $\sum_{n=1}^{N}$ to $ln N(x_n | M_1, \Sigma) = -\frac{1}{2} \sum_{n=1}^{N}$ to $(x_n - M_1)^T \Sigma^{-1}(x_n - M_1) + constant$

Sotting the derivative w.r.t MI, to O, we obtain

$$\mathcal{M}_{i} = \frac{1}{\mathcal{N}_{i}} \sum_{n=1}^{N} t_{n} \mathcal{X}_{n}$$

which is simply the mean of vectors X_n assigned to C_1 .

Similarly we can obtain the result for M_2 as $M_2 = \frac{1}{N_2} \sum_{n=1}^{N} (1-t_n) \times_n$

which again is the mean of vectors X_n assigned to C_2 . Finally, consider the MLE solution for Σ . Pick out the terms in the log likelihood function depending on Z, we have $-\frac{1}{2}\sum_{n=1}^{N} \operatorname{tn} \ln |\Sigma| - \frac{1}{2}\sum_{n=1}^{N} \operatorname{tn} (\mathscr{X}_{n} - \mathcal{M}_{1})^{\mathsf{T}} \Sigma^{\mathsf{T}} (\mathscr{X}_{n} - \mathcal{M}_{1})$ $-\frac{1}{2}\sum_{n=1}^{N}(1-t_n)\ln|\Sigma| -\frac{1}{2}\sum_{n=1}^{N}(1-t_n)(\varkappa_n-\varkappa_1)^{\mathsf{T}}\Sigma^{\mathsf{T}}(\varkappa_n-\varkappa_1)$

$$= -\frac{\sqrt{2}}{2} \ln |\Sigma| - \frac{\sqrt{2}}{2} \operatorname{Tr} \left\{ \Sigma^{-1} S \right\}$$

where we have defined

$$S = \frac{N_1}{N}S_1 + \frac{N_2}{N}S_2$$

$$S_{1} = \frac{1}{N_{1}} \sum_{n \in C_{1}} (X_{n} - M_{1}) (X_{n} - M_{1})^{T}$$

$$S_{2} = \frac{1}{N} \sum_{n \in C_{2}} (x_{n} - M_{2}) (x_{n} - M_{2})^{T}$$

Using the standard result for MLE solution for a Craussian distribution, we see Chapter 4 kinear Models for Classification 64

4.1.3 Discrete features Consider the case of discrete feature value x; For simplicity, assume $x_i \in \{0, 1\}$ and $D-\dim$ vector X $\chi = (\chi_1, \chi_2, \dots, \chi_p)^T$ Here we will make the naive Bayes assumption (the feature values are treated as independent, conditioned on Ck) I.e. $p(X_1, X_2 | C_E) = p(X_1 | C_E) p(X_2 | C_E)$ Thus class - conditional distributions are given by $P(\mathscr{K} | C_{k}) = \prod_{\substack{i=1 \ i \neq i}}^{P} P(X_{i} | C_{k}) = \prod_{\substack{i=1 \ i \neq i}}^{P} \mu_{ki} (I - \mu_{ki}) \qquad (4.81)$ Chapter 4 Linear Models for Classification

which contain D independent parameters for each class.

$$p(x | C_k) = M_k^{\chi} (1 - M_k)^{1-\chi}$$

 $M_k : C_k 2+\xi > + 3\xi_{\delta} + \sigma_1 \chi = 1 \cdot \xi_{\delta}^{\chi}$
 $\chi \in \{0, 1\}$

Substituting into
$$c(4, 63)$$
 ($\alpha_k = ln(pc*(c_k) Pcc_k)$)

$$O_{\mathbf{k}}(\mathbf{x}) = \sum_{i=1}^{n} \int \mathcal{X}_{i} \ln M_{\mathbf{k}i} + (1-\mathbf{x}_{i}) \ln (1-M_{\mathbf{k}i}) + \ln P(C_{\mathbf{k}})$$

which are linear linear functions of xz

4.3 Probabilistic Discriminative Models Finding the parameters of a generalized linear model Grenerative model vs discriminative model (indirect) (direct) Generative model: Fitting class conditional densities PCXICE) priors separately and then applying Baye's Theorem. class and Discriminative model: Maximizing a likelihoot function defined through the conditional distribution PCCELX) Remarks of two approaches

4. 3.1	Fixed	bosis	functions			
⊉ c•) :	vector	of b	rsis functions	ζ¢.,.	- Pm-1 1.	¢. (*) = 1
We M	nake a	fixed	nonlinear t	ansformation	of the	in puts.
The r	esulting.	decision	boundaries	will be	nonlinear	in the
origina	I input	* sp	ace Clinear	in the	feature sp	pace)
We s	hall incl	ude a	fixed basi	s function	transforma	tion $\overline{\mathcal{P}}(x)$.



4.3.2 Logistic regression Two - class classification In section 4.2 (4.57), we saw that under rather general assumptions, the posterior probability of class C, can be written os $p(\zeta | \overline{\Phi}) = \gamma(\overline{\Phi}) = \sigma(w^{T}\overline{\Phi})$ $P(C_1|\overline{\Phi}) = |-P(C_1|\overline{\Phi})$. Here $\sigma(\cdot)$ is the logistic with

sigmoid function and $\overline{\Phi}$ is the feature vector i.e. $\overline{\Phi} = \overline{\Phi}(X)$

For M-Jim feature space, this model has <u>M</u> adjustable parameters (wr).

By contrast, Graussian class conditional densities model using likelihood method needs 2M parameters for maximum mean and M(M+1)/2 parameters vectors for shared covariance Together with class prior this gives a total matrix. of M(M+5)/2 + 1 parameters quadratically

Petermine the parameters of the logistic regression model.
Use maximum likelihood method.
For a dota set
$$d \ Den, tn \ ''$$
 where $tn \in 10, 14$ and $\ Den = \ Den (26n)$
with $n=1, 2, ..., N$, the likelihood function can be written
 $p(t \mid wr) = \prod_{n=1}^{N} y_n^{tn} d \mid -y_n t^{1-tn}$ $tn = 0 \text{ or } 1$
where $t := (t_1, t_2, ..., t_N)^T$ and $y_n = p(C_1 \mid D t) = \sigma(w^T D t)$
 $p(t_1 \mid wr) = y_n^{tn} (1-y_n)^{1-tn}$, $t_n = 0 \text{ or } 1$, its prediction y_n or $1-y_n$
Negative logarithm of the likelihoot which gives the cross entropy error function in the form $E(W) = -\ln P(\#|W) = -\sum_{n=1}^{N} |\tan \ln Y_n + (1-\tan) \ln (1-Y_n)|^2$ where $y_n = \sigma(ca_n)$, $a_n = w_T \overline{\Phi}_n$ with $\overline{\Phi}_n := \overline{\Phi}(w_n)$. The gradient of the error function w.r.t W is given by $\nabla_{wr} E(wr) = \sum_{n=1}^{N} (Y_n - t_n) \Phi_n$ (4.91)ertor

We have used $\frac{4\sigma}{4a} = \sigma(1-\sigma)$

Chapter 4 Linear Models for Classification

From	(4.91),	we	CH	use	a	çe	quential	algorit	hm.	The
weight	vector	W/	`s	uptote	4	in	which	VE,	is	the
n th to	erm in	ር ዓ	91).							

4.3.3	Iterated	reweighted	least	squares		
In the	case of	the line	ear regre	stion mo	tel, MI	E
solution,	on the	assumption	of a	Granssian	noise	model,
leads to	a closet	- form s	olution.			
For logist	tic regres	sion, there	is no	longer	a clos	ed - form
solution.	However t	he error fo	unction E	cwr) is	convex.	Hence
there is	a unique	minimum				

The error function can be minimized by an iterative technique
based on the Newton - Raphson iterative optimization scheme.
(Fletcher 1987) ; Bishop and Nabney 2008)
$$w_{1}^{(new)} = w_{1}^{(OH)} - H^{-1} \nabla E(w_{1})$$

where H is the Hessian matrix $\nabla \nabla E(w_{1})$ w.n.t w.
First, apply Newton - Raphson to the linear regression
model $\nabla_{w_{1}} E(w_{1}) = \sum_{n=1}^{N} (w_{1}^{T} \Phi_{n} - t_{n}) \Phi_{n} = \Phi^{T} \Phi w_{1} - \Phi^{T} \Phi$

$$H = \nabla \nabla Chapter 4 Linear Models for Classification dep of W)$$
 76

where
$$\overline{\Phi}$$
 is the N×M design matrix whose nth row
is given by $\overline{\Phi}_{n}^{T}$. So
 $w_{i}^{(new)} = w_{i}^{(old)} - H^{-1} \nabla E C w_{i}^{(old)}$
 $= w_{i}^{(old)} - (\overline{\Phi}^{T} \overline{\Phi})^{-1} \langle \overline{\Phi}^{T} \overline{\Phi} w_{i}^{(old)} - \overline{\Phi}^{T} t \rangle$
 $= (\overline{\Phi}^{T} \overline{\Phi})^{-1} \langle \overline{\Phi}^{T} t$

is the standard least-squares solution. Since the SSE is the quadratic form of W/, Newton-Raphson formula gives the exact solution in one step

Chapter 4 Linear Models for Classification

Since $Y_n = \sigma(w^T \overline{\phi}(w_n))$,

matrix H is positive definite. Hence the Hessian So erbr function E is a convex function of w and I minimum with $w = w^{(o|d)} - (\overline{\Phi}^T R \overline{\Phi})^T \overline{\Phi}^T (\gamma - t)$ $= (\underline{\Phi}^{\mathsf{T}} R \underline{\Phi})^{\mathsf{T}} \langle \underline{\Phi}^{\mathsf{T}} R \underline{\Phi} w^{\mathsf{cold}} \rangle - \underline{\Phi}^{\mathsf{T}} (\underline{\gamma} - \underline{\mathtt{H}}) \rangle$ $= (\overline{\Phi}^{\mathsf{T}} R \overline{\Phi})^{\mathsf{T}} \overline{\Phi}^{\mathsf{T}} R \mathbf{z}$ (4.99)

where \exists is the N-dimensional vector with $\exists := \Phi W^{(cold)} - R^{((\gamma - t))}$

I least-squared
Recall the MLE solution for
$$W_{1}$$
 of the linear regression
 $W_{ML} = (\overline{\Phi}^T \overline{\Phi})^{-1} \overline{\Phi}^T t$
(4.99) is the form of a set of normal equations for a weighted least - squares problem.

The weighting matrix R is not constant but depends on w.

So we must apply the normal equations iteratively. For this reason, the algorithm is known as IRLS iterative rewighted least squares. 4.3.4 Multiclass logistic regression In section 4.2, we discussed the generative models for multiclass classification. The posterior probabilities are given by a softmax transformation of linear functions of feature variables

$$P(C_{k}|\Phi) = Y_{k}(\Phi) = \frac{\exp(\alpha_{k})}{\Sigma_{j}\exp(\alpha_{j})}$$

where the activations are given by

$$\alpha_{k} = W_{k}^{T} \Phi \qquad \Phi = (\varphi_{0} U_{k}), \dots \varphi_{M-1} U_{k})^{T}$$

There	we	used	MLE	to	determin	ne	se parate	ely the	Graussi	an
class -	con di	itional	densiti	es	and	the	class	priors	and	then
found	the	corre	sponding.	pe	osterior	pro	babilitie	s, H	hereby	
implicitly determining the parameters 1 write										
Here	we	consider	the	use	, र्ज	May	<i>simum</i>	likelihoo	d to	
tetermin	e t	:he po	rameter	5 3	W/k {	र्ठ	this	model	directly	•

 t_n : one hot vector Φ_n : feature vector The lifelihood function is given by $P(T|W_{1},...,W_{k}) = \prod_{n=1}^{N} \prod_{k=1}^{k} P(C_{k}|\Phi)^{tnk} = \prod_{n=1}^{N} \prod_{k=1}^{k} Y_{nk}$ where $Y_{nk} := Y_k(\underline{\mathfrak{P}}_n)$ and T is an NXK matrix of target variables with elements tak $\| p \in H \mid W_1, \dots W_F \| = \prod_{r=1}^{K} p \in C_F \mid \overline{P} \mid H = one hot vector$ ol veight vector = 2 on of target to >+ un find its 84 For some fixed n, Chapter A Linear Models for Classification = 1

Taking the negative logarithm then gives

$$E(w_{1,..}, w_{k}) = -\ln P(T(w_{1,..}, w_{k})) = -\sum_{n=1}^{N} \sum_{k=1}^{k} t_{nk} \ln Y_{nk}$$

is known as the cross-entropy error function which Sor the multiclass classification problem Note that the terivatives of Yr w.r.t all a: $\frac{\partial Y_{k}}{\partial a_{j}} = Y_{k} (I_{kj} - Y_{j}) \qquad \left(Y_{k} = \frac{\exp (a_{k})}{\Sigma_{i} \exp (a_{j})}\right)$ Iki are the elements of identity matrix where

Chapter 4 Linear Models for Classification

take the gradient of the We error function w.r.t now parameter vectors W_i. of the one basis function $\nabla_{W_{l}} \in (W_{l} \dots W_{k}) = \sum (Y_{nj} - t_{nj}) \not \square_{n}$ (4,109)U SI prediction error used Ir tak =1. have where we see the same form arising Note that we for the found gradient was for SSE with linear regression as cross - entropy error for the logistic regression and the model

So we can use this to formulate a sequential algorithm.
In this case, each of the weight vectors is updated
using
$$w_{1}^{(ct+1)} = w_{1}^{(c)} - \eta \nabla E_{n}$$
 (3.22)
Now to find a batch algorithm, we appeal to the

Newton - Raphson update to obtain the corresponding IRLS MK × MK algorithm. The Hessian matrix that comprises blocks of size M×M in which block j, k is given by $\nabla_{W_j} \nabla_{W_k} E(W_{(...,W_k)}) = \sum_{\substack{n=1\\n \geq 1\\ Chapter 4 Linear Models for Classification}}^{N} \nabla_{n_k} \nabla_{$

This	Hessian	matrix	for	the	multic	class	logistic	regression	model
, is	positive	<i>definite</i>	and	50	the	ertor	functi	on again	has
a	unique	minimum.							

4.3.5 Probit reqression two-class classification and the Consider the Stamework of generalized linear models so that $P(t=1|\alpha) = f(\alpha)$ where $a = w^T \Phi$ and f(-) is the activation function. a noisy threshold model. For each input Consider $\Phi_n = \Phi(X_n)$, we evaluate $a_n = W^T \Phi_n$ and then set

the target value according to



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

Chapter 4 Linear Models for Classification

2

3



Chapter 4 Linear Models for Classification

Remark

- It has sigmoidal shape
- The use of a general Gaussian does not change the model
- erf function

$$erf(a) := \frac{2}{\sqrt{2}} \int_{0}^{a} exp(-o^{2}) do$$

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

Chapter 4 Linear Models for Classification

erf function is related to the inverse probit function by $\overline{\Phi}(\alpha) := \frac{1}{2} \int (1 + erf(\frac{\alpha}{\sqrt{2}}))^{1/2}$

The generalized linear model based on an inverse probit activation function is known as probit regression. Remark - The probit model is significantly more sensitive to autliers. - signoit exp(-x) vs inverse probit $exp(-x^2)$ $as x \rightarrow + \infty$

4.4	The	Laplas	ie /	Approxi	mation					
In	Section	4.5	we	will	discuss	the	Bayesian	treatm	ent	of
logis	itic n	egression	n. W	e co	nnot in	tegrate	exactly	over t	the	
para	meter	vector	- W	sin	ce the	: post	erior dist	ribution	is	NO
longe	er Gro	ussian.	So	` t 's	s neces	sary	to introdu	ce so	me	
form	হা	appro.	ximati	on.						
Now	we	intro	dice	the	Laplace	appr	oximation,	that	aims	
to	find	a Go	ussian	app	roximatio	in to	unknown	prob.	densit	ヤ
Jesin	vo be	er a	set	ব	continu	ous	variables.			

Chapter 4 Linear Models for Classification

2: single continuous variable Suppose the distribution pcz) is defined by pcz) = $\frac{1}{Z} fcz$)

where Z is a normalization constant and assumed to be unknown. In the Laplace method, the goal is to find a Growssian approximation Q(2) which is centered on a mode of p(2)

First find a mode p(z), i.e. $z_0 = 0$ $\frac{45(z_0)}{42} = 0$

Note that the logarithm of Gaussian distribution is a quadratic form of variables. Therefore a Taylor expension of h f(z) centered on the mode zo is given by

$$lnf(z) \simeq lnf(z_{0}) - \frac{1}{2}A(z-z_{0})^{2}$$

where

$$A = -\frac{4^{2}}{4^{2}} \ln f(2)$$

Chapter 4 Linear Models for Classification

Taking the exponential we obtain $f(z) \simeq f(z_0) \exp \left(-\frac{A}{2}(z_0 - z_0)^2 \right)$

We can then obtain a normalized distribution Q(z) so that $Q(z) = \left(\frac{A}{2\pi}\right)^{1/2} \exp\left(1 - \frac{A}{2}(z - z_0)^2\right)$

Note that it will only be well defined if its precision A > 0 (20 must be local maximum or 5''(20) < 0)

₹: M-dim vector

the Laplace method to approximate $P(\mathcal{X}) = f(\mathcal{X})/2$ Extend At a stationary point zo, Vf(z) will vanish. Expanding around this stationary point to we have $l_n f(z) \simeq l_n f(z) - \frac{1}{2}(z - z_0)^T A (z - \overline{z}_0)$ MXM Hessian matrix A is defined by where. **^**

$$A = -\nabla \nabla \ln f(x)|_{x=x_0}$$

Taking exponential us obtain $f(z) \simeq f(z_0) \exp(-\frac{1}{2}(z-z_0)^T A(z-z_0))$

Thus

$$q(2) = \frac{|A|^{2}}{(2\pi)^{M/2}} \exp\left(-\frac{1}{2}(2-2)^{T}A(2-2)\right) = N(2|2_{0}, A^{T})$$

where IAI denotes the determinant of A. As before, this Graussian will be well defined if A is positive definite.

Remark

- Need to find a mode to and evaluate Hessian matrix.
- In practice, a mode will be sound by running some

form of numerical optimization algorithm

- Limitations of multimodal case
- Normalization constant 2 does not need to be known.

As well as approximating the distribution
$$p(z_{2})$$
, we can
obtain an approximation to z
 $z = \int f(z_{2}) dz \simeq f(z_{2}) \int exp \{-\frac{1}{2}(z_{2}-z_{2})^{T} \wedge (z_{2}-z_{2}) \{dz \}$
 $= f(z_{2}) \frac{(2\pi)^{M_{2}}}{|A|^{M_{2}}}$

4.5	Baye	esian	Logisti	ic Rev	gression)				
The	evalu	nation	र्ज र	the pos	terior	distribut	tion	over \	M wou	19
requi	ire r	ormaliza	tion	of the	protu	it of	0	prior	distribut	zion
and	0	likelihood	Sun	ction.	Note 1	that f	the	i Felihoo	d funct	ion
compr	rises	a pro	duct .	of log	istic s	tigmoit	(by	our	assumpti	ion j
j.e.	PC#	(w) =	Π _n γ,	tn $(1 - \gamma_n)$	I-tn)	$\gamma_n = c$	5 (WT	₽_n).		
Evalue	ation	of the	e pr	edictive	distri	bution	is s	imilarly	intrac	table.
Here	we	consider	the	applic	ation	of the	e La	place	approximo	ction
to	the	problem	रु	Boyesi	an lo	gistic	regre	ession		

4.5.1 Laplace approximation We need the evaluation of the second derivatives of the log posterior (finding the Hessian matrix) we seek a Graussian representation (approximation) Beauce for the posterior distribution, we introduce a Graussian prior. $pcwr) = NCwr (M_0, S_0)$ mo, So are fixed hyperparameters. where

The posterior distribution over
$$w$$
 is given by
prior likelihood
 $P(wr1\#) \circ C P(wr) P(\#1wr)$
Generican product of signal
where $\# := (t_1, ..., t_N)^T$. Taking the log of both sides.
 $l_n P(wr1\#) = -\frac{1}{2}(wr - m_b)^T s_0^{-1}(wr - m_b)$
 $+ \sum_{n=1}^{N} i_1 t_n l_n y_n + (1 - t_n) l_n(1 - y_n) i_1 + constant$
where $y_n = e^{-(wr - m_b)}$.

To obtain a Gaussian approximation to the posterior
distribution, we first maximize the posterior distribution
to give the MAP (moximum a posterior) solution
$$W_{MAP}$$

defining the mean (mode) of Gaussian. The covariance is
then given by
 $S_{N}^{-1} = -\nabla \nabla \ln PCW(1 \#) = S_{0}^{-1} + \sum_{n=1}^{N} Y_{n} (1 - Y_{n}) \oint_{n} \oint_{n}^{T}$
The Gaussian approximation to the posterior distribution

$$q(w) = \mathcal{N}(w| | w|_{MAP}, S_{\mathcal{N}})$$

Chapter 4 Linear Models for Classification

4.5.2 Predictive distribution There remains the take of marginalizing w.r.t q cm) to make prediction. Let $\overline{Q} = \overline{Q}(x)$ be the seature vector. The predictive distribution for C1 is obtained by marginalizing. w.r.t p(w/1#), which is itself approximated by a Graussian distribution qCW) so that $PCC(\overline{\Phi}, \overline{R}) = \int PCC(\overline{\Phi}, w) PCW(\overline{R}) 4w \sim \int ccw_T \overline{\Phi}) qcw) 4w$ i.e. $PCC_2(\overline{Q}, \underline{H}) = I - PCC_1(\overline{Q}, \underline{H})$

Let $S(\cdot)$ be the Dirac delta function. Then we have $\sigma(w^T \overline{\Phi}) = \int S(a - w^T \overline{\Phi}) \sigma(a) da$

From this

$$\int \sigma cw T \Phi (w) dw = \int \sigma ca product da$$

= $E_a [\sigma]$

where

$$P(\alpha) = \int S(\alpha - w^T \overline{\Phi}) Q(w) dw$$

new prob. distribution

The Dirac tetta imposes a linear constraint on wr.
So
$$p(\alpha)$$
 forms a marginal distribution from the joint
distribution $q(wr)$ by integrating out all directions
orthogonal to $\overline{\Phi}$. It follows that $p(\alpha)$ is Gaussian.
 $M_{\alpha} = E[\alpha] = \int p(\alpha) \alpha d\alpha = \int \int S(\alpha - wT\overline{\Phi}) q(wr) dwr \alpha d\alpha$

$$= \int \int S(a - w^{T} \Phi) a daq(w) dw$$

$$= \int w^{T} \Phi q(w) dw$$

$$q(w) = N(w) | w_{MAP}, S_{N})$$

$$= W_{MAP} \Phi$$
Similarly,

$$\sigma_{\alpha}^{2} = \operatorname{Var}[\alpha] = \int p(\alpha) \left\{ a^{2} - \operatorname{E}[\alpha] \right\}^{2} d\alpha$$

= $\int q(w) \left\{ (w^{T} \overline{\Phi})^{2} - (m_{w}^{T} \overline{\Phi})^{2} \right\} \left\{ dw = \overline{\Phi}^{T} S_{w} \overline{\Phi}$
Ne have used $q(w) = N(w) \left\{ W_{MAP}, S_{w} \right\}.$

Thus, the variational approximation to the predictive distribution becomes (4.151)

$$P(\zeta_1 \#) \simeq \int \sigma(\alpha) p(\alpha) d\alpha = \int \sigma(\alpha) N(\alpha) M(\alpha, \sigma_{\alpha}^2) d\alpha$$

This integral cannot be analytically. So we approximate

$$\sigma(\alpha)$$
 by $\overline{P}(n\alpha)$ with suitable value $n(\alpha) n^2 = \frac{\pi}{8}$
The advantage of using an inverse profit function is that
the below integral (convolution) can be expressed analytically
in terms of another inverse profit function.
 $\int \overline{P}(n\alpha) N(\alpha | M, \sigma^2) d\alpha = \overline{P}\left(\frac{M}{(n^{-2} + \sigma^2)^{\frac{1}{2}}}\right)$
(Spiegel halter and Lauritzen 1990; Mackay 1992 b; Barber
and Bishop, 1998a)

We apply the approximation $\sigma(\alpha) \cong \overline{\Phi}(n\alpha)$ and it leads to the following approximation $\int \overline{\sigma}(\alpha) N(\alpha) \mu, \sigma^2) d\alpha \cong \sigma(K(\sigma^2)\mu)$

where we defined

$$K(\sigma^{2}) = (1 + \pi \sigma^{2}/8)$$

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

 $P(C(|\overline{\Phi}, \#) \simeq \sigma(k(\overline{\sigma_a}) \mu_a)$