

# Weekly Overview Slides of Statistical Machine Learning CSE 575, Fall 2023

Moses A. Boudourides<sup>1</sup>

SPA and SCAI  
Arizona State University

<sup>1</sup> [Moses.Boudourides@asu.edu](mailto:Moses.Boudourides@asu.edu)

## Week 6

*Exercises on Linear Basis Function Models, Discriminant  
Functions, and Probabilistic Generative Models*

February 16, 2023

# Exercise 1

## The logistic sigmoid function in terms of tanh

Show that the logistic sigmoid function  $\sigma(\alpha) = 1/(1 + \exp(-\alpha))$  satisfies the equation  $2\sigma(2\alpha) - 1 = \tanh$  and find the corresponding form to the general linear combination of logistic sigmoid function  $y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \sigma\left(\frac{x - \mu_j}{s}\right)$ .

## Solution

$$\begin{aligned} 2\sigma(2\alpha) - 1 &= \frac{2}{1 + e^{-2\alpha}} - 1 = \frac{2}{1 + e^{-2\alpha}} - \frac{1 + e^{-2\alpha}}{1 + e^{-2\alpha}} \\ &= \frac{1 - e^{-2\alpha}}{1 + e^{-2\alpha}} = \frac{e^\alpha - e^{-\alpha}}{e^\alpha + e^{-\alpha}} = \tanh(\alpha). \end{aligned}$$

Thus, taking  $\alpha_j = (x - \mu_j)/(2s)$ , we get

$$\begin{aligned} y(\mathbf{x}, \mathbf{w}) &= w_0 + \sum_{j=1}^M w_j \sigma(2\alpha_j) = w_0 + \sum_{j=1}^M \frac{w_j}{2} (2\sigma(2\alpha_j) - 1 + 1) \\ &= u_0 + \sum_{j=1}^M u_j \tanh(\alpha_j), \end{aligned}$$

where  $u_j = w_j/2$ , for  $j = 1, \dots, M$ , and  $u_0 = w_0 + \sum_{j=1}^M w_j/2$ .

## Exercise 2

### Orthogonal Projections in Linear Algebra

If the columns of a matrix  $\Phi$  are independent, show that the matrix  $\Phi(\Phi^T\Phi)^{-1}\Phi^T$  takes any vector  $\mathbf{v}$  and projects it onto the space spanned by the columns of  $\Phi$ .

### Solution

Since the columns of  $\Phi$  are independent, one can show by contradiction that  $\Phi^T\Phi$  is invertible. Because, if  $\Phi^T\Phi$  was not invertible, then there would exist a vector  $\mathbf{x} \neq \mathbf{0}$  such that  $\Phi^T\Phi\mathbf{x} = \mathbf{0}$ . So,  $\mathbf{x}^T(\Phi^T\Phi\mathbf{x}) = 0$ , i.e.,  $(\Phi\mathbf{x})^T(\Phi\mathbf{x}) = 0$ , which would mean that  $\|\Phi\mathbf{x}\| = 0$ , i.e.,  $\Phi\mathbf{x} = \mathbf{0}$ , and, thus, the columns of  $\Phi$  could not be independent, which is a contradiction. Moreover, the invertibility of  $\Phi^T\Phi$  implies that  $(\Phi^T\Phi)^{-1} = \Phi^{-1}\Phi^{-T}$ . Next, let us consider the matrix

$$\Psi = \Phi(\Phi^T\Phi)^{-1}\Phi^T.$$

Then it is easy to see that both the following two conditions hold

$$\Psi^T = \Psi$$

$$\Psi^2 = \Psi.$$

Therefore,  $\Psi$  is a projection matrix on the range of  $\Phi$ , which is what we want.

## Exercise 3

### A Linear Model

Consider a linear model of the form

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i, \quad (1)$$

together with a sum-of-squares error function of the form

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \quad (2)$$

Now suppose that Gaussian noise  $\epsilon_i$  with zero mean and variance  $\sigma^2$  is added independently to each of the input variables  $x_i$ . By making use of  $\mathbb{E}[\epsilon_i] = 0$  and  $\mathbb{E}[\epsilon_i \epsilon_j] = \delta_{ij} \sigma^2$ , show that minimizing  $E_D$  averaged over the noise distribution is equivalent to minimizing the sum-of-squares error for noise-free input variables with the addition of a weight-decay regularization term, in which the bias parameter  $w_0$  is omitted from the regularizer.

Setting

$$\begin{aligned}\tilde{y}_n &= w_0 + \sum_{i=1}^D w_i (x_{ni} + \epsilon_{ni}) \\ &= y_n + \sum_{i=1}^D w_i \epsilon_{ni},\end{aligned}$$

where  $y_n = y(x_n, \mathbf{w})$  and  $\epsilon_{ni} \sim \mathcal{N}(0, \sigma^2)$ , after using (1). Thus, as in (2), we define

$$\begin{aligned}\tilde{E} &= \frac{1}{2} \sum_{n=1}^N \{\tilde{y}_n - t_n\}^2 = \frac{1}{2} \sum_{n=1}^N \{\tilde{y}_n^2 - 2\tilde{y}_n t_n + t_n^2\} \\ &= \frac{1}{2} \sum_{n=1}^N \left\{ y_n^2 + 2y_n \sum_{i=1}^D w_i \epsilon_{ni} + \left( \sum_{i=1}^D w_i \epsilon_{ni} \right)^2 - 2t_n y_n - 2t_n \sum_{i=1}^D w_i \epsilon_{ni} + t_n^2 \right\}.\end{aligned}$$

Taking the expectation of  $\tilde{E}$  under the distribution of  $\epsilon_{ni}$ , the second and the fifth term disappear, since  $\mathbb{E}[\epsilon_{ni}] = 0$ , while the third term becomes

$$\mathbb{E} \left[ \left( \sum_{n=1}^D w_i \epsilon_{ni} \right)^2 \right] = \sum_{i=1}^D w_i^2 \sigma^2,$$

since the  $\epsilon_{ni}$  are all independent with variance  $\sigma^2$ . From this and (2), we obtain

$$\mathbb{E}[\tilde{E}] = E_D + \frac{1}{2} \sum_{i=1}^D w_i^2 \sigma^2.$$

## Exercise 4

### Lagrange Multipliers

Using the technique of Lagrange multipliers, show that the minimization of the regularized error function (3.29) is equivalent to minimizing the unregularized sum-of-squares error (3.12) subject to the constraint

$$\sum_{j=1}^M |w_j|^q \leq \eta.$$

### Solution

Apparently, the above constraint can be written as  $\frac{1}{2} \left( \sum_{j=1}^M |w_j|^q - \eta \right) \leq 0$ , which combined with (3.12) yields the Lagrange function

$$L(\mathbf{w}, \lambda) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \left( \sum_{j=1}^M |w_j|^q - \eta \right),$$

which is identical to (3.29) in the dependence on  $\mathbf{w}$ . Now, let us choose a specific value of  $\lambda > 0$  to minimize (3.29). Denoting the resulting value of  $\mathbf{w}$  as  $\mathbf{w}^*(\lambda)$ , and using the KKT condition  $\lambda g(\mathbf{w}) = 0$ , we get the following value of  $\eta$

$$\eta = \sum_{j=1}^M |w_j^*(\lambda)|^q.$$

## Exercise 5

### A Linear Basis Function Regression Gaussian Model

Consider a linear basis function regression model for a multivariate target variable  $\mathbf{t}$  having a Gaussian distribution of the form

$$p(\mathbf{t}|\mathbf{W}, \Sigma) = \mathcal{N}(\mathbf{t}|\mathbf{y}(\mathbf{x}, \mathbf{W}), \Sigma),$$

where

$$\mathbf{y}(\mathbf{x}, \mathbf{W}) = \mathbf{W}^T \phi(\mathbf{x}),$$

together with a training data set comprising input basis vectors  $\phi(\mathbf{x}_n)$  and corresponding target vectors  $\mathbf{t}_n$ , with  $n = 1, \dots, N$ . Show that the maximum likelihood solution  $\mathbf{W}_{\text{ML}}$  for the parameter matrix  $\mathbf{W}$  has the property that each column is given by an expression of the form (3.15), which was the solution for an isotropic noise distribution. Note that this is independent of the covariance matrix  $\Sigma$ . Show that the maximum likelihood solution for  $\Sigma$  is given by

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{\text{ML}}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{\text{ML}}^T \phi(\mathbf{x}_n))^T.$$

## Solution of Exercise 5

The corresponding log-likelihood function is

$$\log L(\mathbf{W}, \Sigma) = -\frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{n=1}^N \left( \mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n) \right)^T \Sigma^{-1} \left( \mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n) \right).$$

Setting the  $\mathbf{W}$  derivative equal to 0 yields

$$0 = - \sum_{n=1}^N \Sigma^{-1} \left( \mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n)^T.$$

Multiplying through by  $\Sigma$  and introducing the design matrix  $\Phi$  and the target data matrix  $\mathbf{T}$  we have

$$\Phi^T \Phi \mathbf{W} = \Phi^T \mathbf{T},$$

which solved for  $\Phi$  gives

$$\mathbf{W}_{\text{ML}} = \left( \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{T},$$

where the operator  $\left( \Phi^T \Phi \right)^{-1} \Phi^T = \Phi^\dagger$  is the *Moore–Penrose pseudo-inverse* of the matrix  $\Phi$ . Finally, the maximum likelihood solution for  $\Sigma$  is found as in section 3.1 (by appealing to the standard result from Chapter 2).



## Exercise 6

### Convex Hulls and Linear Separability

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

$$\mathbf{x} = \sum_n \alpha_n \mathbf{x}_n,$$

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

## Solution of Exercise 6

If the convex hulls of  $\{\mathbf{x}_n\}$  and  $\{\mathbf{y}_n\}$  intersect, there exist a point  $\mathbf{z}$  such that  $\mathbf{z} = \sum_n \alpha_n \mathbf{x}_n = \sum_n \beta_n \mathbf{y}_n$ . Hence, using that the coefficients of a convex hull sum up to 1, we would have

$$\begin{aligned}\hat{\mathbf{w}}^T \mathbf{z} + w_0 &= \hat{\mathbf{w}}^T \left( \sum_n \alpha_n \mathbf{x}_n \right) + w_0 = \left( \sum_n \alpha_n \hat{\mathbf{w}}^T \mathbf{x}_n \right) + \left( \sum_n \alpha_n \right) w_0 \\ &= \sum_n \alpha_n (\hat{\mathbf{w}}^T \mathbf{x}_n + w_0).\end{aligned}$$

Now, if  $\{\mathbf{x}_n\}$  and  $\{\mathbf{y}_n\}$  were linearly separable, we would have  $\hat{\mathbf{w}}^T \mathbf{x}_n + w_0 > 0$ , for all  $\mathbf{x}_n$ , and  $\hat{\mathbf{w}}^T \mathbf{y}_n + w_0 < 0$ , for all  $\mathbf{y}_n$ . Clearly, the first of the latter conditions would imply (using the previous expression) that  $\hat{\mathbf{w}}^T \mathbf{z} + w_0 > 0$ . Similarly, from the corresponding properties of  $\{\mathbf{y}_n\}$ , we would get that  $\hat{\mathbf{w}}^T \mathbf{z} + w_0 < 0$ , which is a contradiction and, therefore, the required result follows.

## Exercise 7

### Least Squares for Classification

Consider the minimization of a sum-of-squares error function  $E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr}\{(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})\}$  (\*), and suppose that all of the target vectors in the training set satisfy a linear constraint

$$\mathbf{a}^T \mathbf{t}_n + b = 0,$$

where  $\mathbf{t}_n$  corresponds to the  $n$ -th row of the matrix  $\mathbf{T}$  in (\*). Show that, as a consequence of this constraint, the elements of the model prediction  $\mathbf{y}(\mathbf{x})$  given by the least-squares solution  $\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T (\tilde{\mathbf{X}}^\dagger)^T \tilde{\mathbf{x}}$  (where  $\tilde{\mathbf{X}}^\dagger$  is the pseudo-inverse of the matrix  $\tilde{\mathbf{X}}$ ) also satisfy this constraint, so that

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

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

## Solution of Exercise 7

Remember that in the  $N$ -dimensional classification case, since  $\tilde{\mathbf{W}} = [\mathbf{w}_0, \mathbf{W}^T]^T$  and  $\tilde{\mathbf{X}} = [\mathbf{1}, \mathbf{X}^T]^T$ , we get the following explicit dependence of  $E_D(\tilde{\mathbf{W}})$  on  $\mathbf{w}_0$ :

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \{ (\mathbf{XW} + \mathbf{1w}_0^T - \mathbf{T})^T (\mathbf{XW} + \mathbf{1w}_0^T - \mathbf{T}) \}.$$

The above expression is a matrix quadratic w.r.t.  $\mathbf{w}_0$  and its  $\mathbf{w}_0$  derivative can be easily found to be:

$$\frac{\partial}{\partial \mathbf{w}_0} E_D(\tilde{\mathbf{W}}) = 2N\mathbf{w}_0 + 2(\mathbf{XW} - \mathbf{T})^T \mathbf{1}.$$

Thus, setting the derivative equal to 0, we get

$$\mathbf{w}_0 = -\frac{1}{N}(\mathbf{XW} - \mathbf{T})^T \mathbf{1} = \bar{\mathbf{t}} - \mathbf{W}^T \bar{\mathbf{x}},$$

where

$$\bar{\mathbf{t}} = \frac{1}{N} \mathbf{T}^T \mathbf{1} \text{ and } \bar{\mathbf{x}} = \frac{1}{N} \mathbf{X}^T \mathbf{1}.$$

Substituting the above critical value of  $\mathbf{w}_0$  in the expression for  $E_D(\tilde{\mathbf{W}})$ , we obtain

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \{ (\mathbf{XW} + \bar{\mathbf{T}} - \bar{\mathbf{X}}\mathbf{W} - \mathbf{T})^T (\mathbf{XW} + \bar{\mathbf{T}} - \bar{\mathbf{X}}\mathbf{W} - \mathbf{T}) \},$$

where we have denoted

$$\bar{\mathbf{T}} = \mathbf{1}\bar{\mathbf{t}}^T \text{ and } \bar{\mathbf{X}} = \mathbf{1}\bar{\mathbf{x}}^T.$$

## Solution of Exercise 7 (cont.)

Next, denoting two new matrices

$$\hat{\mathbf{X}} = \mathbf{X} - \bar{\mathbf{X}} \text{ and } \hat{\mathbf{T}} = \mathbf{T} - \bar{\mathbf{T}},$$

setting the  $\mathbf{W}$  derivative of  $E_D(\tilde{\mathbf{W}})$  equal to 0 yields

$$\mathbf{W} = (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \hat{\mathbf{T}} = \hat{\mathbf{X}}^\dagger \hat{\mathbf{T}},$$

where  $\hat{\mathbf{X}}^\dagger = (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T$  is the *Moore–Penrose pseudo-inverse* of  $\hat{\mathbf{X}}$ .

Finally, consider the prediction for a new input vector  $\mathbf{x}^*$

$$\begin{aligned} \mathbf{y}(\mathbf{x}^*) &= \mathbf{W}^T \mathbf{x}^* + \mathbf{w}_0 = \mathbf{W}^T \mathbf{x}^* + \bar{\mathbf{t}} - \mathbf{W}^T \bar{\mathbf{x}} \\ &= \bar{\mathbf{t}} + \hat{\mathbf{T}}^T (\hat{\mathbf{X}}^\dagger)^T (\mathbf{x}^* - \bar{\mathbf{x}}). \end{aligned}$$

If we know that  $\mathbf{a}^T \mathbf{t}_n + b = 0$  holds, for some  $\mathbf{a}$  and  $b$ , we get

$$\mathbf{a}^T \bar{\mathbf{t}} = \frac{1}{N} \mathbf{a}^T \mathbf{T}^T \mathbf{1} = -b.$$

Therefore, applying the condition  $\mathbf{a}^T \mathbf{t}_n + b = 0$  to the previously derived expression for  $\mathbf{x}^*$ , we obtain

$$\mathbf{a}^T \mathbf{y}(\mathbf{x}^*) = \mathbf{a}^T \bar{\mathbf{t}} + \mathbf{a}^T \hat{\mathbf{T}}^T (\hat{\mathbf{X}}^\dagger)^T (\mathbf{x}^* - \bar{\mathbf{x}}) = \mathbf{a}^T \bar{\mathbf{t}} = -b,$$

since  $\mathbf{a}^T \hat{\mathbf{T}}^T = \mathbf{a}^T (\mathbf{T} - \bar{\mathbf{T}})^T = b(\mathbf{1} - \mathbf{1})^T = \mathbf{0}^T$ .

## Exercise 8

### Maximization of the Class Separation Criterion with Lagrange Multiplier

Show that maximization with respect to  $\mathbf{w}$  of the class separation criterion given in Bishop's PRML book as  $m_2 - m_1 = \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1)$ , using a Lagrange multiplier to enforce the constraint  $\mathbf{w}^T \mathbf{w} = 1$ , leads to the result that  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ .

### Solution

Using a Lagrange multiplier to enforce the constraint  $\mathbf{w}^T \mathbf{w} = 1$  means that we need to maximize the following objective function

$$L(\lambda, \mathbf{w}) = \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1) + \lambda(\mathbf{w}^T \mathbf{w} - 1).$$

Clearly, the  $\mathbf{w}$  derivative (gradient) of this Lagrange function is

$$\frac{\partial}{\partial \mathbf{w}} L(\lambda, \mathbf{w}) = \mathbf{m}_2 - \mathbf{m}_1 + 2\lambda \mathbf{w},$$

which set equal to 0 gives

$$\mathbf{w} = -\frac{1}{2\lambda}(\mathbf{m}_2 - \mathbf{m}_1) \propto (\mathbf{m}_2 - \mathbf{m}_1).$$

## Exercise 9

### The Alternative Form of the Fisher Criterion

By making use of (4.20), (4.23), and (4.24), show that the Fisher criterion (4.25) can be written in the form (4.26), where the previous parentheses correspond to numbers of equations in Bishop's PRML book.

### Solution

First, we expand (4.25) using (4.22), (4.23) and (4.24):

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\|\mathbf{w}^T(m_2 - m_1)\|^2}{\sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - m_1)^2 + \sum_{n \in C_2} (\mathbf{w}^T \mathbf{x}_n - m_2)^2}.$$

Next, on the one hand, denoting

$$\mathbf{S}_B = (m_2 - m_1)(m_2 - m_1)^T,$$

the numerator becomes

$$\text{numerator} = (\mathbf{w}^T(m_2 - m_1))(\mathbf{w}^T(m_2 - m_1))^T = \mathbf{w}^T \mathbf{S}_B \mathbf{w}.$$

On the other hand, denoting

$$\mathbf{S}_w = \sum_{n \in C_1} (\mathbf{x}_n - m_1)(\mathbf{x}_n - m_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - m_2)(\mathbf{x}_n - m_2)^T,$$

the denominator becomes

$$\begin{aligned} \text{denominator} &= \sum_{n \in C_1} (\mathbf{w}^T(\mathbf{x}_n - m_1))^2 + \sum_{n \in C_2} (\mathbf{w}^T(\mathbf{x}_n - m_2))^2 \\ &= \mathbf{w}^T \mathbf{S}_{w_1} \mathbf{w} + \mathbf{w}^T \mathbf{S}_{w_2} \mathbf{w} = \mathbf{w}^T \mathbf{S}_w \mathbf{w}. \end{aligned}$$

Thus, the wanted result.

## Exercise 10

### The Logistic Sigmoid Function

Show that the logistic sigmoid function  $\sigma(\alpha) = 1/(1 + \exp(-\alpha))$  satisfies the property  $\sigma(-\alpha) = 1 - \sigma(\alpha)$  and that its inverse is given by  $\sigma^{-1}(y) = \ln\{y/(1 - y)\}$ .

### Solution

By definition of the logistic sigmoid function, we get

$$\begin{aligned}\sigma(\alpha) + \sigma(-\alpha) &= \frac{1}{1 + \exp(-\alpha)} + \frac{1}{1 + \exp(\alpha)} \\ &= \text{DO THE ALGEBRA} = 1.\end{aligned}$$

Next, exchanging the dependent and independent variables in  $y = \sigma(\alpha)$ , we obtain the inverse of the logistic sigmoid function

$$\alpha = \frac{1}{1 + \exp(-y)}$$

and after rearranging it as

$$\exp(-y) = \frac{1 - \alpha}{\alpha},$$

if we take the logarithm in both sides, we get

$$y = \ln\left\{\frac{\alpha}{1 - \alpha}\right\},$$

which is what we wanted (writing  $\sigma^{-1}(y)$  instead of  $y$  and  $y$  instead of  $\alpha$ ).



# Exercise 11

## The Two-Class Generative Model with Gaussian Densities

Using (4.57) and (4.58), derive the result (4.65) for the posterior class probability in the two-class generative model with Gaussian densities, and verify the results (4.66) and (4.67) for the parameters  $\mathbf{w}$  and  $w_0$  (where all parentheses correspond to numbers of equations in Bishop's PRML book).

### Solution

According to (4.58) and (4.64), we can write

$$\begin{aligned}\alpha &= \ln \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)} = \ln p(\mathbf{x}|C_1) - \ln p(\mathbf{x}|C_2) + \ln \frac{p(C_1)}{p(C_2)} \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} \\ &= \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)} \\ &= \mathbf{w}^T \mathbf{x} + w_0.\end{aligned}$$

Notice that, in the last step, we have rearranged the terms, and have set

$$\begin{aligned}\mathbf{w} &= \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \\ w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}.\end{aligned}$$

Thus, since  $p(C_1|\mathbf{x}) = \sigma(\alpha)$  (as stated in (4.57)), we get  $p(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$ , as required.

## Exercise 12

### Maximum–Likelihood Solution to a Generative Classification Model

Consider a generative classification model for  $K$  classes defined by prior class probabilities  $p(C_k) = \pi_k$  and general class–conditional densities  $p(\phi|C_k)$  where  $\phi$  is the input feature vector. Suppose we are given a training data set  $\{\phi_n, t_n\}$ , where  $n = 1, \dots, N$ , and  $t_n$  is a binary target vector of length  $K$  that uses the 1–of– $K$  coding scheme, so that it has components  $t_{nj} = I_{jk}$ , if pattern  $n$  is from class  $C_k$ . Assuming that the data points are drawn independently from this model, show that the maximum–likelihood solution for the prior probabilities is given by

$$\pi_k = \frac{N_k}{N},$$

where  $N_k$  is the number of data points assigned to class  $C_k$ .

The likelihood function is given by

$$p(\{\phi_n, \mathbf{t}_n\}|\{\pi_k\}) = \prod_{n=1}^N \prod_{k=1}^K \{p(\phi_n|C_k)\pi_k\}^{t_{nk}}$$

and taking the logarithm we obtain

$$\ln p(\{\phi_n, \mathbf{t}_n\}|\{\pi_k\}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \{\ln p(\phi_n|C_k) + \ln \pi_k\} \propto \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln \pi_k.$$

Since there is a constraint on  $\pi_k$ , we need to add a Lagrange multiplier to the above expression

$$L = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln \pi_k + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right).$$

Clearly, the  $\pi_k$  derivative of the Lagrangian function is

$$\frac{\partial}{\partial \pi_k} L = \sum_{n=1}^N \frac{t_{nk}}{\pi_k} + \lambda$$

and, setting the derivative equal to 0, we obtain

$$\pi_k = -\frac{\sum_{n=1}^N t_{nk}}{\lambda} = -\frac{N_k}{\lambda}.$$

Summing both sides of the last expression with regard to  $k$ , we get

$$1 = -\frac{\sum_{n=1}^N N_k}{\lambda} = -\frac{N}{\lambda},$$

i.e.,  $\lambda = -N$ , which implies that  $\pi_k = N_k/N$ .