

The Elements of Statistical Learning - Chapter 4 Exercises

Exercise 4.1

Show how to solve the generalized eigenvalue problem $\max a^T \mathbf{B} a$ subject to $a^T \mathbf{W} a = 1$ by transforming to a standard eigenvalue problem.

Solution

Take the eigendecomposition $\mathbf{W} = \mathbf{U} \mathbf{D} \mathbf{U}^T$ where \mathbf{U} is a $p \times p$ orthogonal matrix. Let $\mathbf{M}^* = \mathbf{M} \mathbf{D}^{-\frac{1}{2}}$ and $a^* = \mathbf{D}^{\frac{1}{2}} \mathbf{U}^T a$. The covariance matrix of \mathbf{M}^* can be obtained from that of \mathbf{M} by conjugating:

$$\mathbf{B}^* = \mathbf{U}^T \mathbf{D}^{-\frac{1}{2}} \mathbf{B} \mathbf{D}^{\frac{1}{2}} \mathbf{U},$$

so the constrained optimisation problem is equivalent to $\max a^{*T} \mathbf{B}^* a^*$ subject to $a^{*T} a^* = 1$.

Consider the eigendecomposition $\mathbf{B}^* = \mathbf{V}^* \mathbf{D}_B \mathbf{V}^{*T}$ where the diagonal entries of \mathbf{D}_B are $\hat{\mu}_1, \dots, \hat{\mu}_p$ (they are all positive). Then $a^{*T} \mathbf{B}^* a^* = \|\mathbf{D}_B^{\frac{1}{2}} \mathbf{V}^{*T} a^*\|^2$. Write a^* as a linear combination of the eigenvectors: $a^* = \sum_{i=1}^p \hat{\lambda}_i v_i^*$. This transforms the optimisation problem into the following: $\max \sum \hat{\lambda}_i^2$ subject to $\sum \hat{\lambda}_i^2 = 1$. This is a standard Lagrange multiplier problem with solution $a^* = v_1^*$.

Transforming back to the original coordinates gives

$$a = \mathbf{U} \mathbf{D}^{-\frac{1}{2}} a^* = \mathbf{U} \mathbf{D}^{-\frac{1}{2}} v_1^* = v_1$$

as required.

Exercise 4.2

Suppose we have features $x \in \mathbb{R}^p$, a two-class response, with class sizes N_1, N_2 , and the targets coded as $-N/N_1, N/N_2$.

(a) Show that the LDA rule classifies to class 2 if

$$x^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) > \frac{1}{2} (\hat{\mu}_2 + \hat{\mu}_1)^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) - \log(N_2/N_1),$$

and class 1 otherwise.

(b) Consider minimization of the least squares criterion

$$\sum_{i=1}^N (y_i - \beta_0 - x_i^T \beta)^2.$$

Show that the solution $\hat{\beta}$ satisfies

$$\left[(N-2) \hat{\Sigma} + N \hat{\Sigma}_B \right] \beta = N(\hat{\mu}_2 - \hat{\mu}_1)$$

(after simplification), where $\hat{\Sigma}_B = \frac{N_1 N_2}{N^2} (\hat{\mu}_2 - \hat{\mu}_1)(\hat{\mu}_2 - \hat{\mu}_1)^T$.

(c) Hence show that $\hat{\Sigma}_B \beta$ is in the direction $(\hat{\mu}_2 - \hat{\mu}_1)$ and thus

$$\hat{\beta} \propto \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1).$$

Therefore the least-squares regression coefficient is identical to the LDA coefficient, up to a scalar multiple.

(d) Show that this result holds for any (distinct) coding of the two classes.

(e) Find the solution $\hat{\beta}_0$ (up to the same scalar multiple as in (c)), and hence the predicted value $\hat{f}(x) = \hat{\beta}_0 + x^T \hat{\beta}$. Consider the following rule: classify to class 2 if $\hat{f}(x) > 0$ and class 1 otherwise. Show this is not the same as the LDA rule unless the classes have equal numbers of observations.

Solution

(a) The discriminant functions for LDA are

$$\delta_k(x) = x \hat{\Sigma}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\Sigma}^{-1} \hat{\mu}_k + \ln(\pi_k).$$

Estimating μ_k by $\hat{\mu}_k$, Σ by $\hat{\Sigma}$, and π_k by N_k/N , we classify to class 2 if and only if

$$\begin{aligned} \hat{\delta}_2(x) &> \hat{\delta}_1(x) \\ \Leftrightarrow x \hat{\Sigma}^{-1} \hat{\mu}_2 - \frac{1}{2} \hat{\Sigma}^{-1} \hat{\mu}_2 + \ln(N_2/N) &> x \hat{\Sigma}^{-1} \hat{\mu}_1 - \frac{1}{2} \hat{\Sigma}^{-1} \hat{\mu}_1 + \ln(N_1/N) \\ \Leftrightarrow x^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) &> \frac{1}{2} (\hat{\mu}_2 - \hat{\mu}_1)^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) + \ln(N_1/N_2) \\ \Leftrightarrow x^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) &> \frac{1}{2} (\hat{\mu}_2 + \hat{\mu}_1)^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) - \log(N_2/N_1) \end{aligned}$$

and to class 1 otherwise.

(b) Let $\tilde{\mathbf{X}}$ denote the input data matrix augmented with the intercept column. We know that $\hat{\beta}$ satisfies $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \hat{\beta} = \tilde{\mathbf{X}}^T \mathbf{y}$.

We wish to relate $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \hat{\beta}$ and $\mathbf{X}^T \mathbf{X} \beta$. The (j, k) entry of $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is $\sum_{i=1}^N x_{ij} x_{ik}$. Taking $k=0$, the j th entry of the 0th column is

$$\sum_i x_{ij} x_{i0} = \sum_{i=1}^N x_{ij} + \sum_{i=2}^N x_{ij} = N_1 \hat{\mu}_{1j} + N_2 \hat{\mu}_{2j}.$$

So $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is a $(p+1) \times (p+1)$ -matrix with $(0, 0)$ entry N , $(1, 0), \dots, (p, 0)$ entries taken up by $N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2$, $(0, 1), \dots, (0, p)$ entries taken up by $N_1 \hat{\mu}_1^T + N_2 \hat{\mu}_2^T$, and bottom right $p \times p$ submatrix equal to $\mathbf{X}^T \mathbf{X}$. This implies that $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \hat{\beta}$ is equal to the block matrix

$$\begin{pmatrix} N \hat{\beta}_0 + (N_1 \hat{\mu}_1^T + N_2 \hat{\mu}_2^T) \beta \\ \beta_0 (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2) + \mathbf{X}^T \mathbf{X} \beta \end{pmatrix},$$

where the top expression is a scalar and the bottom is a $p \times 1$ vector.

We have already seen that this is equal to $\tilde{\mathbf{X}}^T \mathbf{y}$. The j th entry of $\tilde{\mathbf{X}}^T \mathbf{y}$ is

$$x_j^T \mathbf{y} = \frac{N}{N_2} \left(\sum_{i \in \mathcal{C}_2} x_{ij} \right) - \frac{N}{N_1} \left(\sum_{i \in \mathcal{C}_1} x_{ij} \right)$$

so $\tilde{\mathbf{X}}^T \mathbf{y}$ is $N(\hat{\mu}_2 - \hat{\mu}_1)$ augmented with a zero in the 0th place. Equating the expressions in the 0th place gives

$$\hat{\beta}_0 = -\frac{1}{N} (N_1 \hat{\mu}_1^T + N_2 \hat{\mu}_2^T) \beta.$$

Equating the remaining terms and substituting this in gives

$$\left[\mathbf{X}^T \mathbf{X} - \frac{1}{N} (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2)^T (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2) \right] \beta = N(\hat{\mu}_2 - \hat{\mu}_1).$$

It remains to show that the expression in square brackets equals $(N-2) \hat{\Sigma} + N \hat{\Sigma}_B$.

Indeed,

$$\begin{aligned} (N-2) \hat{\Sigma} &= \sum_{k=1}^2 \sum_{i \in \mathcal{C}_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T \\ &= \sum_{k=1}^2 \sum_{i \in \mathcal{C}_k} (x_i x_i^T - \hat{\mu}_k x_i^T - x_i \hat{\mu}_k^T + \hat{\mu}_k \hat{\mu}_k^T) \\ &= \sum_{k=1}^2 \left[\left(\sum_{i \in \mathcal{C}_k} x_i x_i^T \right) - N_k \hat{\mu}_k \hat{\mu}_k^T - N_k \hat{\mu}_k \hat{\mu}_k^T + N_k \hat{\mu}_k \hat{\mu}_k^T \right] \\ &= \mathbf{X}^T \mathbf{X} - \sum_{k=1}^2 N_k \hat{\mu}_k \hat{\mu}_k^T. \end{aligned}$$

So,

$$\begin{aligned} (N-2) \hat{\Sigma} + N \hat{\Sigma}_B &= \mathbf{X}^T \mathbf{X} - N_1 \hat{\mu}_1 \hat{\mu}_1^T - N_2 \hat{\mu}_2 \hat{\mu}_2^T + \frac{N_1 N_2}{N} (\hat{\mu}_2 - \hat{\mu}_1)(\hat{\mu}_2 - \hat{\mu}_1)^T \\ &= \mathbf{X}^T \mathbf{X} + \left(\frac{N_1 N_2}{N} - N_1 \right) \hat{\mu}_1 \hat{\mu}_1^T + \left(\frac{N_1 N_2}{N} - N_2 \right) \hat{\mu}_2 \hat{\mu}_2^T - \frac{N_1 N_2}{N} (\hat{\mu}_1 \hat{\mu}_2^T + \hat{\mu}_2 \hat{\mu}_1^T) \\ &= \mathbf{X}^T \mathbf{X} - \frac{1}{N} (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2)(N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2)^T. \end{aligned}$$

as required.

(c) By definition,

$$\hat{\Sigma}_B \beta = \frac{N_1 N_2}{N^2} (\hat{\mu}_2 - \hat{\mu}_1)(\hat{\mu}_2 - \hat{\mu}_1)^T \beta = \left[\frac{N_1 N_2}{N^2} (\hat{\mu}_2 - \hat{\mu}_1)^T \beta \right] (\hat{\mu}_2 - \hat{\mu}_1).$$

So

$$\hat{\Sigma} \hat{\beta} = \frac{N}{N-2} \left[1 - \frac{N_1 N_2}{N^2} (\hat{\mu}_2 - \hat{\mu}_1)^T \beta \right] (\hat{\mu}_2 - \hat{\mu}_1)$$

and thus $\hat{\beta} \propto \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1)$.

(d) Changing the coding makes little difference to our solution to part (b). Suppose that the targets are coded as c_1 and c_2 . Then

$$\mathbf{X}^T \mathbf{y} = \begin{pmatrix} c_1 N_1 + c_2 N_2 \\ c_1 N_1 \hat{\mu}_1 + c_2 N_2 \hat{\mu}_2 \end{pmatrix}.$$

Taking this through our solution we end up with

$$\left[(N-2) \hat{\Sigma} + N \hat{\Sigma}_B \right] \beta + \frac{c_1 N_1 + c_2 N_2}{N} (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2) = c_1 N_1 \hat{\mu}_1 + c_2 N_2 \hat{\mu}_2.$$

Rearranging this gives

$$\left[(N-2) \hat{\Sigma} + N \hat{\Sigma}_B \right] \beta = (c_2 - c_1) \frac{N_1 N_2}{N} (\hat{\mu}_2 - \hat{\mu}_1).$$

This implies that $\hat{\beta} \propto \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1)$ as before.

(e) By part (b), $\hat{\beta}_0 = -\frac{1}{N} (N_1 \hat{\mu}_1^T + N_2 \hat{\mu}_2^T) \beta$, so

$$\hat{\beta}_0 \propto -\frac{1}{N} (N_1 \hat{\mu}_1^T + N_2 \hat{\mu}_2^T) \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1)$$

with the same constant of proportionality as $\beta \propto \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1)$. Therefore

$$\hat{f}(x) \propto \left(x - \frac{1}{N} (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2) \right)^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1).$$

So we classify to class 2 iff

$$x^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) > \frac{1}{N} (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2)^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1).$$

Since $N = N_1 + N_2$, this is the same as the LDA rule if $N_1 = N_2$. If they are not equal, the rules will be different in general.

Exercise 4.3

Suppose we transform the original predictors \mathbf{X} to $\hat{\mathbf{Y}}$ via linear regression. In detail, let $\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{X} \hat{\mathbf{B}}$, where \mathbf{Y} is the indicator response matrix. Similarly for any input $x \in \mathbb{R}^p$, we get a transformed vector $\hat{y} = \hat{\mathbf{B}}^T x \in \mathbb{R}^K$. Show that LDA using $\hat{\mathbf{Y}}$ is identical to LDA in the original space.

Solution

Let $\hat{\mu}_k$, $\hat{\mu}_k$, and $\hat{\Sigma}$ denote the parameter estimates for the original predictors. Clearly \hat{x} is unchanged under the transformation. The new estimates of class means are

$$\frac{1}{N_k} \sum_{i \in \mathcal{C}_k} \hat{y}_i = \frac{1}{N_k} \sum_{i \in \mathcal{C}_k} \hat{\mathbf{B}}^T x_i = \hat{\mathbf{B}}^T \hat{\mu}_k.$$

Since

$$\hat{\Sigma} = \frac{1}{N-K} \left(\mathbf{X}^T \mathbf{X} - \sum_{k=1}^K N_k \hat{\mu}_k \hat{\mu}_k^T \right),$$

the transformed covariance estimate is

$$\frac{1}{N-K} \left(\hat{\mathbf{Y}}^T \hat{\mathbf{Y}} - \sum_{k=1}^K N_k \hat{\mathbf{B}}^T \hat{\mu}_k \hat{\mu}_k^T \hat{\mathbf{B}} \right) = \hat{\mathbf{B}}^T \hat{\Sigma} \hat{\mathbf{B}}.$$

Therefore the new linear discriminant functions are

$$\begin{aligned} \delta_k(\hat{y}) &= (x^T \hat{\mathbf{B}}) \left(\hat{\mathbf{B}}^{-1} \hat{\Sigma}^{-1} (\hat{\mathbf{B}}^T)^{-1} \right) (\hat{\mathbf{B}}^T \hat{\mu}_k) - \frac{1}{2} (\hat{\mu}_k^T \hat{\mathbf{B}}) \left(\hat{\mathbf{B}}^{-1} \hat{\Sigma}^{-1} (\hat{\mathbf{B}}^T)^{-1} \right) (\hat{\mathbf{B}}^T \hat{\mu}_k) + \ln(\hat{\pi}_k) \\ &= x^T \hat{\Sigma} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \hat{\Sigma} \hat{\mu}_k + \ln(\hat{\pi}_k) \\ &= \delta_k(x), \end{aligned}$$

as required.

Exercise 4.4

Consider the multilogit model with K classes (4.17). Let β be the $(p+1)(K-1)$ -vector consisting of all the coefficients. Define a suitably enlarged version of the input vector x to accommodate this vectorized coefficient matrix. Derive the Newton-Raphson algorithm for maximizing the multinomial log-likelihood, and describe how you would implement this algorithm.

Solution

Given a length p input vector x (we assume the intercept 1 is already included), let \tilde{x} be the $(p+1)(K-1) \times (K-1)$ block matrix with $K-1$ x 's along the diagonal and 0s elsewhere. Then $\tilde{x}^T \beta$ is the length $(K-1)N$ vector with k th entry $\beta_k^T x$ (again we include the intercept term in β_k). From this the vector of probabilities $p_k(x; \theta)$ can easily be calculated.

We now derive the score equations and Newton-Raphson algorithm. Let $y_i = (y_{i0}, \dots, y_{iK})^T$ denote the i th row of the indicator response matrix \mathbf{Y} . The multinomial log-likelihood is

$$\begin{aligned} l(\beta) &= \sum_{i=1}^N \sum_{k=1}^K y_{ik} \ln(p_k(x_i; \theta)) \\ &= \sum_{i=1}^N \left[\sum_{k=1}^{K-1} \left[y_{ik} \left(\beta_k^T x_i - \ln(1 + \sum_{l=1}^{K-1} \exp(\beta_l^T x_i)) \right) \right] - \left(1 - \sum_{l=1}^{K-1} y_{il} \right) \ln(1 + \sum_{l=1}^{K-1} \exp(\beta_l^T x_i)) \right] \\ &= \sum_{i=1}^N \left[\sum_{k=1}^{K-1} \left(y_{ik} \beta_k^T x_i \right) - \ln(1 + \sum_{l=1}^{K-1} \exp(\beta_l^T x_i)) \right]. \end{aligned}$$

Differentiating this with respect to one of the β_k gives

$$\begin{aligned} \frac{\partial l(\beta)}{\partial \beta_k} &= \sum_{i=1}^N \left[y_{ik} x_i - \frac{\exp(\beta_k^T x)}{1 + \sum_{l=1}^{K-1} \exp(\beta_l^T x_i)} x_i \right] \\ &= \sum_{i=1}^N (y_{ik} - p_k(x_i; \theta)) x_i \end{aligned}$$

and differentiating again yields

$$\frac{\partial^2 l(\beta)}{\partial \beta_k \partial \beta_l^T} = \begin{cases} -\sum_{i=1}^N p_k(x_i; \theta) p_l(x_i; \theta) x_i x_i^T & \text{if } k \neq l \\ -\sum_{i=1}^N p_k(x_i; \theta) (1 - p_k(x_i; \theta)) x_i x_i^T & \text{if } k = l \end{cases}$$

This gives us enough to implement Newton-Raphson: concatenate the $\frac{\partial l(\beta)}{\partial \beta_k}$ and $\frac{\partial^2 l(\beta)}{\partial \beta_k \partial \beta_l^T}$ into a single vector and block matrix $\frac{\partial l(\beta)}{\partial \beta}$ and $\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T}$, respectively. Then set

$$\beta^{\text{new}} = \beta^{\text{old}} - \left(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial l(\beta)}{\partial \beta}.$$

We will now describe how to use the enlarged vector \tilde{x} to put everything in terms of matrix products. The idea is to vectorise natural matrices like \mathbf{Y} in a similar fashion to β .

Let $\tilde{\mathbf{X}}$ denote the $(k-1)N \times (p+1)(K-1)$ -matrix obtained by stacking $\tilde{x}_1^T, \dots, \tilde{x}_N^T$ on top of each other, or equivalently obtained from \mathbf{X} by replacing each x_i^T with its tile counterpart. Let $\tilde{\mathbf{y}}$ be the length $(k-1)N$ vector obtained by stacking the observations y_i on top of each other in order. Similarly, let $\tilde{\mathbf{p}}$ denote the length $(K-1)$ vector with entries $p_1(x_i; \theta), \dots, p_{K-1}(x_i; \theta)$ and let $\tilde{\mathbf{p}}$ be the length $(k-1)N$ vector obtained by stacking the $\tilde{\mathbf{p}}_i$. Then

$$\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^N \tilde{x}_i (y_i - \tilde{\mathbf{p}}_i) = \tilde{\mathbf{X}}^T (\tilde{\mathbf{y}} - \tilde{\mathbf{p}})$$

For $i \in \{1, \dots, N\}$, let \mathbf{W}_i denote the $(K-1) \times (K-1)$ matrix with (k, l) entry $p_k(x_i; \theta) p_l(x_i; \theta)$ if $k \neq l$ and $p_k(x_i; \theta)(1 - p_k(x_i; \theta))$ if $k = l$. Then

$$\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^N \tilde{x}_i \mathbf{W}_i \tilde{x}_i^T = \tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}},$$

where \mathbf{W} is the $(K-1)N \times (K-1)N$ block matrix with $\mathbf{W}_1, \dots, \mathbf{W}_N$ along the diagonal.

With this in hand, the Newton-Raphson algorithm can be expressed as an iterated reweighted least squares procedure as in the $K=2$ case:

$$\beta^{\text{new}} = (\tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{z}}, \quad \text{where} \quad \tilde{\mathbf{z}} = \tilde{\mathbf{X}} \beta^{\text{old}} + \mathbf{W}^{-1} (\tilde{\mathbf{y}} - \tilde{\mathbf{p}})$$

This can be implemented directly, updating \mathbf{W} and $\tilde{\mathbf{p}}$ at each step, or as the solution to the weighted least squares problem $\min_{\tilde{\beta}} (\tilde{\mathbf{z}} - \tilde{\mathbf{X}} \tilde{\beta})^T \mathbf{W} (\tilde{\mathbf{z}} - \tilde{\mathbf{X}} \tilde{\beta})$.

Exercise 4.5

Consider a two-class logistic regression problem with $x \in \mathbb{R}$. Characterize the maximum-likelihood estimates of the slope and intercept parameter if the sample x_i for the two classes are separated by a point $x_0 \in \mathbb{R}$. Generalize this result to $\{\mathbf{a}\} x \in \mathbb{R}^p$ (see Figure 4.16), and (b) more than two classes.

Solution

In all these situations the MLEs are undefined. In the first case the log-likelihood is

$$l(\beta) = \sum_{i=1}^N [y_i \ln(p(x_i; \theta)) + (1 - y_i) \ln(1 - p(x_i; \theta))].$$

Since $p(x; \theta) \in [0, 1]$, this is negative. If the two classes are separated then we can approach 0 with suitable choices of θ , but θ will diverge to infinity.

To be more precise, suppose that $y_i = 1$ if $x_i < x_0$ and $y_i = 0$ if $x_i > x_0$. Take $\theta = (\beta_0, \beta) = (-rx_0, r)$ for $r > 0$. Then $\hat{\beta}_0 + \hat{\beta}x > 0$ if $x > x_0$ and $\hat{\beta}_0 + \hat{\beta}x > 0$ if $x < x_0$. So as $r \rightarrow \infty$,

$$\hat{\beta}_0 + \hat{\beta}x_i \rightarrow \begin{cases} -\infty & \text{if } y_i = 1 \\ \infty & \text{if } y_i = 0 \end{cases} \implies p(x_i; \theta) \rightarrow \begin{cases} 0 & \text{if } y_i = 1 \\ 1 & \text{if } y_i = 0 \end{cases}$$

and $l(\beta) \rightarrow 0$. Clearly this path isn't unique, for example we could take a different separating point x_0 .

(a) If $x \in \mathbb{R}^p$ for $p > 1$ the situation is very similar. Suppose that the two classes are separated by a hyperplane $\tilde{\beta}_0 + \tilde{\beta}^T x = 0$ with $y_i = 1$ if $\hat{\beta}_0 + \hat{\beta}^T x < 0$ and $y_i = 0$ if $\hat{\beta}_0 + \hat{\beta}^T x > 0$. Then $\theta = (r\hat{\beta}_0, r\hat{\beta})$ diverges to infinity as $r \rightarrow \infty$ but $l(\beta)$ converges to its supremum zero\$.

(b) There are a few ways to conceive of generalising to $K > 2$ classes, but the point is that the log-likelihood has supremum 0 which can't be attained with finite values of θ . For a simple formulation, suppose that for $k = 1, \dots, K-1$, the k th class can be separated from all other classes by a hyperplane $\hat{\beta}_{0k} + \hat{\beta}_k x = 0$. More precisely, suppose that $g_i = k$ if $\hat{\beta}_{0k} + \hat{\beta}_k x_i > 0$ and $g_i \neq k$ if $\hat{\beta}_{0k} + \hat{\beta}_k x_i < 0$. Let θ have k th component $(r\hat{\beta}_{0k}, r\hat{\beta}_k)$. Then the log-likelihood

$$l(\beta) = \sum_{i=1}^N \sum_{k=1}^K y_{ik} \ln(p_k(x_i; \theta))$$

converges to 0 as $r \rightarrow \infty$ but θ diverges to infinity.

Note that this means that logistic regression fails for arguably the simplest type of classification problem. This motivates separating hyperplanes.

Exercise 4.6

Suppose we have N points x_i in \mathbb{R}^p in general position, with class labels $y_i \in \{-1, 1\}$. Prove that the perceptron learning algorithm converges to a separating hyperplane in a finite number of steps:

(a) Denote a hyperplane by $f(x) = \beta_0^T + \beta_0 = 0$, or in more compact notation $\beta^T x^* = 0$, where $x^* = (x, 1)$ and $\beta = (\beta_1, \beta_0)$. Let $z_i = x_i^* \|x_i^*\|$. Show that separability implies the existence of a β_{sep} such that $y_i \beta_{\text{sep}}^T z_i \geq 1 \forall i$