

# Statistical learning homework 2

Dai Yuehao (1800010660@pku.edu.cn)

November 2, 2020

## 1 ESL 3.16

For best subset selection we select  $M$  coefficients into our model, thus we have

$$\hat{y}' = \sum_{i=1}^p \hat{\beta}_i I_i x_i,$$

here  $I_i = 1$  if  $\hat{\beta}_i$  is chosen into the model,  $I_i = 0$  if not. Now the residual is

$$r = \|y - \hat{y}'\|_2^2 = \|y\|_2^2 - \sum_{i=1}^p \hat{\beta}_i^2 I_i^2 \|x_i\|_2^2 = \|y\|_2^2 - \sum_{i=1}^p \hat{\beta}_i^2 I_i,$$

hence the best selection of coefficients should be the  $M$  largest ones, hence the formula is

$$\hat{\beta}'_j = \hat{\beta}_j \cdot I[\text{rank}(|\hat{\beta}_j| \leq M)].$$

For ridge regression we have the close form solution

$$\hat{\beta}' = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T y = (1 + \lambda) \mathbf{X}^T y,$$

since

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y = \mathbf{X}^T y$$

we have the formula

$$\hat{\beta}'_j = \frac{\hat{\beta}_j}{1 + \lambda}.$$

For LASSO we have the corresponding optimization problem

$$\min f(\beta) = \frac{1}{2} \|y - \mathbf{X}\beta\|_2^2 + \lambda \sum_{i=1}^p \beta_i,$$

this is a convex problem hence the optimality condition is

$$\mathbf{X}^T y - \mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T y - \beta \in \partial \lambda \|\beta\|_1, \quad (1.0.1)$$

since  $\mathbf{X}^T y = \hat{\beta}$  and if  $\hat{\beta}_j > \lambda$  then  $\partial|\hat{\beta}_j| = 1$  and we should shrink  $\hat{\beta}_j$  to  $\hat{\beta}_j - \lambda$  in order to satisfy (1.0.1). Similar operation when  $\hat{\beta}_j < -\lambda$ . If  $0 \leq \hat{\beta}_j \leq \lambda$  we can only shrink it to 0 and the same operation when  $-\lambda \leq \hat{\beta}_j \leq 0$ . Thereby we obtain the formula

$$\hat{\beta}'_j = \text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda)_+.$$

## 2 ESL 3.30

Let

$$\tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda\alpha}\mathbf{I} \end{pmatrix}, \quad \tilde{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}$$

then we have

$$\|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\alpha\|\boldsymbol{\beta}\|_2^2,$$

hence the elastic-net optimization problem can be turned into a lasso problem written as

$$\min_{\boldsymbol{\beta}} \quad \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2 + \lambda(1 - \alpha)\|\boldsymbol{\beta}\|_1.$$

## 3 ADMM for group LASSO

Let  $\boldsymbol{\gamma} = (\boldsymbol{\gamma}_1^T, \dots, \boldsymbol{\gamma}_G^T)^T$ , then we have the equivalent optimization problem

$$\min_{\boldsymbol{\beta}, \boldsymbol{\gamma}} \quad \frac{1}{2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{g=1}^G \|\boldsymbol{\gamma}_g\|_2 \quad \text{s.t.} \quad \boldsymbol{\beta}_g = \boldsymbol{\gamma}_g, \quad g = 1, \dots, G$$

then the augmented Lagrange function is

$$L_\rho(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\alpha}) = \frac{1}{2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{g=1}^G \|\boldsymbol{\gamma}_g\|_2 + \boldsymbol{\alpha}^T(\boldsymbol{\beta} - \boldsymbol{\gamma}) + \frac{\rho}{2}\|\boldsymbol{\beta} - \boldsymbol{\gamma}\|_2^2,$$

then we can renew the coefficients

$$\boldsymbol{\beta}^{(k+1)} \leftarrow (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{y} + \rho \boldsymbol{\gamma}^{(k)} - \boldsymbol{\alpha}^{(k)}),$$

then for given  $\boldsymbol{\gamma}_g$  it is a convex optimization problem

$$\min_{\boldsymbol{\gamma}_g} \quad \lambda \|\boldsymbol{\gamma}_g\|_2 - \boldsymbol{\alpha}^T \boldsymbol{\gamma}_g + \frac{\rho}{2} \|\boldsymbol{\beta}_g - \boldsymbol{\gamma}_g\|_2^2,$$

the optimality condition is

$$\rho(\boldsymbol{\gamma}_g - \boldsymbol{\beta}_g) - \boldsymbol{\alpha} \in -\lambda \partial \|\boldsymbol{\gamma}_g\|_2,$$

then in terms of this we renew  $\boldsymbol{\gamma}_g^{(k)}$  from  $g = 1$  to  $g = G$ , at last we renew

$$\boldsymbol{\alpha}^{(k+1)} \leftarrow \boldsymbol{\alpha}^{(k)} + \rho (\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\gamma}^{(k+1)}).$$

## 4 Normal bound

By the assumption we have

$$P\{Z \geq t\}e^{t^2/(2\sigma^2)} = \int_t^\infty \frac{1}{\sqrt{2\pi}\sigma} e^{-(x^2-t^2)/(2\sigma^2)} dx = \int_0^\infty \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} e^{-xt/\sigma^2} dx,$$

hence

$$P\{Z \geq t\}e^{t^2/(2\sigma^2)} \leq \int_0^\infty \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} dx = \frac{1}{2}.$$

On the other hand for all  $\varepsilon > 0$  there exists  $N > 0$  such that

$$\int_0^N \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} dx > \frac{\sqrt{1-\varepsilon}}{2}$$

holds and there exists  $t_0 > 0$  such that for all  $t < t_0$  and  $x \leq N$

$$e^{-xt/\sigma^2} > \sqrt{1-\varepsilon}$$

holds, hence

$$\int_0^\infty \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} e^{-xt/\sigma^2} dx > \int_0^N \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/(2\sigma^2)} e^{-xt/\sigma^2} dx > \frac{1-\varepsilon}{2},$$

thereby we have

$$\sup_{t>0} \left\{ P\{Z \geq t\} e^{t^2/(2\sigma^2)} \right\} = \frac{1}{2}.$$

## 5 ESL 4.2

### 5.1 a

For class  $k = 1, 2$  and by the Gaussian assumption we can calculate the Bayes' discriminant function

$$\delta_k(x) = -\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) - \frac{p}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_k|) + \ln(\pi_k),$$

by the assumption we let

$$\Sigma_1 = \Sigma_2 = \hat{\Sigma}, \quad \mu_k = \hat{\mu}_k, \quad \pi_k = \frac{N_k}{N}, \quad k = 1, 2,$$

thus the LDA rule classifies  $x$  to class 2 if  $\delta_2(x) - \delta_1(x) > 0$ . Note that on can write

$$\delta_k(x) = -\frac{1}{2} \left( x^T \hat{\Sigma}^{-1} x - x^T \hat{\Sigma}^{-1} \hat{\mu}_k - \hat{\mu}_k^T \hat{\Sigma}^{-1} x + \hat{\mu}_k^T \hat{\Sigma}^{-1} \hat{\mu}_k \right) + \log \left( \frac{N_k}{N} \right),$$

hence

$$\delta_2(x) - \delta_1(x) = x^T \hat{\Sigma}^{-1} (\hat{\mu}_2 - \hat{\mu}_1) - \frac{1}{2} \left( \hat{\mu}_2^T \hat{\Sigma}^{-1} \hat{\mu}_2 - \hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1 \right) - \log \left( \frac{N_1}{N} \right) + \log \left( \frac{N_2}{N} \right)$$

which is exactly our result.

### 5.2 b

Assume that  $x_1, \dots, x_{N_1}$  are in class 1 the rest are in class 2, then we directly calculate that

$$X^T X = \begin{pmatrix} N & \sum_{i=1}^N x_i^T \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i x_i^T \end{pmatrix}, \quad X^T y = \begin{pmatrix} 0 \\ -N\hat{\mu}_1 + N\hat{\mu}_2 \end{pmatrix},$$

also we have

$$\hat{\Sigma} = \frac{1}{N-2} \left[ \sum_{i=1}^n x_i x_i^T - N_1 \hat{\mu}_1 \hat{\mu}_1^T - N_2 \hat{\mu}_2 \hat{\mu}_2^T \right],$$

thus

$$\sum_{i=1}^n x_i x_i^T = (N-2)\hat{\Sigma} + N_1\hat{\mu}_1\hat{\mu}_1^T + N_2\hat{\mu}_2\hat{\mu}_2^T.$$

Since

$$X^T X \begin{pmatrix} \beta_0 \\ \beta \end{pmatrix} = \begin{pmatrix} 0 \\ -N\hat{\mu}_1 + N\hat{\mu}_2 \end{pmatrix}$$

we have

$$N\beta_0 + \sum_{i=1}^N x_i^T \beta = 0 \quad \Rightarrow \quad \beta_0 = \left( -\frac{N_1}{N}\hat{\mu}_1^T - \frac{N_2}{N}\hat{\mu}_2^T \right) \beta, \quad (5.2.1)$$

since  $\sum_{i=1}^N x_i = N_1\hat{\mu}_1 + N_2\hat{\mu}_2$  we have

$$\begin{aligned} \sum_{i=1}^N \beta_0 x_i + \sum_{i=1}^N x_i x_i^T \beta &= -N\hat{\mu}_1 + N\hat{\mu}_2 \\ \Rightarrow \left[ (N-2)\hat{\Sigma} + \frac{N_1 N_2}{N}(\hat{\mu}_1 - \hat{\mu}_2)^T(\hat{\mu}_1 - \hat{\mu}_2)\hat{\Sigma}_B \right] \beta &= N(\hat{\mu}_1 - \hat{\mu}_2), \end{aligned}$$

thus we obtain

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

### 5.3 c

Note that

$$\hat{\Sigma}_B \beta = (\hat{\mu}_1 - \hat{\mu}_2)(\hat{\mu}_1 - \hat{\mu}_2)^T \beta \propto \hat{\mu}_1 - \hat{\mu}_2,$$

thus

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

### 5.4 d

Assume that we place the label  $y_i$  for class  $i$ , then we also calculate that

$$X^T y = \begin{pmatrix} N_1 y_1 + N_2 y_2 \\ y_1 N_1 \hat{\mu}_1 + y_2 N_2 \hat{\mu}_2 \end{pmatrix},$$

then we have

$$N\beta_0 + \sum_{i=1}^N x_i^T \beta = N_1 y_1 + N_2 y_2,$$

hence

$$\beta_0 = \frac{N_1}{N} y_1 + \frac{N_2}{N} y_2 - \left( \frac{N_1}{N} \hat{\mu}_1^T + \frac{N_2}{N} \hat{\mu}_2^T \right) \beta,$$

since

$$\sum_{i=1}^N \beta_0 x_i + \sum_{i=1}^N x_i x_i^T \beta = y_1 N_1 \hat{\mu}_1 + y_2 N_2 \hat{\mu}_2$$

we obtain

$$\left[ (N-2)\hat{\Sigma} + \frac{N_1 N_2}{N}\hat{\Sigma}_B \right] \beta = y_1 N_1 \hat{\mu}_1 + y_2 N_2 \hat{\mu}_2 - (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2) \left( \frac{N_1}{N} y_1 + \frac{N_2}{N} y_2 \right),$$

then we see that

$$y_1 N_1 \hat{\mu}_1 + y_2 N_2 \hat{\mu}_2 - (N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2) \left( \frac{N_1}{N} y_1 + \frac{N_2}{N} y_2 \right) = \frac{N_1 N_2}{N} (y_1 - y_2) (\hat{\mu}_1 - \hat{\mu}_2)$$

then we conclude our proof.

## 5.5 e

From (5.2.1) we have

$$\hat{\beta}_0 = -\frac{1}{N} \sum_{i=1}^N x_i^T \beta.$$

Next we have

$$\hat{f} = \left( x - \frac{1}{N} \sum_{i=1}^N x_i^T \right)^T \beta \propto x^T \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_2) - \frac{1}{N} \sum_{i=1}^N x_i^T \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_2),$$

when  $N_1 = N_2$  the right hand side of above can be written as

$$x^T \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_2) - \frac{1}{2} (\hat{\mu}_1 + \hat{\mu}_2)^T \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_2)$$

which from ESL 4.5(a) we know it is the decision boundary of LDA. When  $N_1 \neq N_2$  they are obviously different.

## 6 ESL 4.3

We can write the simplified discriminant function of  $x$  as

$$\delta_k^{(x)}(x) = x^T \hat{\Sigma}_x^{-1} \hat{\mu}_k^{(x)} - \frac{1}{2} \hat{\mu}_k^{(x)T} \hat{\Sigma}_x^{-1} \hat{\mu}_k^{(x)} + \log \pi_k,$$

and the discriminant function of  $y$  as

$$\delta_k^{(y)}(y) = y^T \hat{\Sigma}_y^{-1} \hat{\mu}_k^{(y)} - \frac{1}{2} \hat{\mu}_k^{(y)T} \hat{\Sigma}_y^{-1} \hat{\mu}_k^{(y)} + \log \pi_k,$$

since  $y = \hat{\mathbf{B}}^T x$  and  $\hat{\mathbf{Y}} = \mathbf{X} \hat{\mathbf{B}}$  we have

$$\hat{\mu}_k^{(y)} = \hat{\mathbf{B}}^T \hat{\mu}_k^{(x)}, \quad \hat{\Sigma}_y = \frac{1}{N - K} \sum_{k=1}^K \sum_{g_i=k} (y_i - \hat{\mu}_i^{(y)})(y_i - \hat{\mu}_i^{(y)})^T = \hat{\mathbf{B}}^T \hat{\Sigma}_x \hat{\mathbf{B}}, \quad (6.0.1)$$

our goal is to verify that  $\delta_k^{(x)}(x) - \delta_l^{(x)}(x) = 0$  if and only if  $\delta_k^{(y)}(x \hat{\mathbf{B}}^T) - \delta_l^{(y)}(x \hat{\mathbf{B}}^T) = 0$ , if  $K \leq p$  we know that  $\hat{\Sigma}_y$  is nonsingular, and from (6.0.1) we need to verify

$$\hat{\mathbf{B}}(\hat{\mathbf{B}}^T \hat{\Sigma}_x \hat{\mathbf{B}})^{-1} \hat{\mathbf{B}}^T (\hat{\mu}_k^{(x)} - \hat{\mu}_l^{(x)}) = \hat{\Sigma}_x^{-1} (\hat{\mu}_k^{(x)} - \hat{\mu}_l^{(x)}).$$

Since  $\mathbf{Y}$  is an indicator response matrix we have

$$\hat{\mu}_k^{(x)} = \frac{1}{N_k} \sum_{g_i=k} x_i = \frac{1}{N_k} \mathbf{X}^T y_k, \quad (6.0.2)$$

thus

$$\hat{\Sigma}_x = \frac{1}{N-K} \left[ \sum_{i=1}^N x_i x_i^T - \sum_{k=1}^K N_k \hat{\mu}_k^{(x)} \hat{\mu}_k^{(x)T} \right] = \frac{1}{N-K} (\mathbf{X}^T \mathbf{X} - \mathbf{P})$$

where

$$\mathbf{P} = \sum_{k=1}^K \frac{1}{N_k} \mathbf{X}^T y_i y_i^T \mathbf{X},$$

hence

$$\hat{\Sigma}_y = \frac{1}{N-K} \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X} - \mathbf{P}) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y},$$

let  $\mathbf{H} = \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$  then

$$(\hat{\mathbf{B}}^T \hat{\Sigma}_x \hat{\mathbf{B}})^{-1} = (N-K)(\mathbf{I} - \mathbf{H}_1)^{-1} \mathbf{H}^{-1}$$

where  $\mathbf{H}_1 = \mathbf{H}^{-1} \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{P} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ , hence from (6.0.2) if we have

$$\begin{aligned} & \hat{\Sigma}_x \hat{\mathbf{B}} (\hat{\mathbf{B}}^T \hat{\Sigma}_x \hat{\mathbf{B}})^{-1} \hat{\mathbf{B}}^T \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{Y} - \mathbf{P} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}) (\mathbf{I} - \mathbf{H}_1)^{-1} \mathbf{H}^{-1} \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \mathbf{X}^T \mathbf{Y}, \end{aligned} \tag{6.0.3}$$

then

$$\hat{\Sigma}_x \hat{\mathbf{B}} (\hat{\mathbf{B}}^T \hat{\Sigma}_x \hat{\mathbf{B}})^{-1} \hat{\mathbf{B}}^T \mathbf{X}^T y_k = \hat{\Sigma}_x \hat{\mathbf{B}} (\hat{\mathbf{B}}^T \hat{\Sigma}_x \hat{\mathbf{B}})^{-1} \hat{\mathbf{B}}^T N_k \hat{\mu}_k^x = \mathbf{X}^T y_k = N_k \hat{\mu}_k^x$$

and we finish the proof. Equation (6.0.3) holds if and only if

$$\mathbf{P} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{X}^T \mathbf{Y} \mathbf{H}_1, \tag{6.0.4}$$

a sufficient condition is  $N_{k_1} = N_{k_2}$  for all  $k_1, k_2$ , under this condition we have

$$\mathbf{P} = \frac{1}{N_1} \mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}$$

and (6.0.4) holds. General cases have not been solved yet.

If  $K > p$  then  $\hat{\Sigma}_y$  is singular and it seems that LDA with  $y$  can not be done. However if we consider that  $\hat{\mathbf{Y}}$  is restricted in a  $p$  dimensional subspace and we do the LDA in that subspace, then we augment  $\mathbf{X}$  to

$$\tilde{\mathbf{X}} = (\mathbf{X} \quad \mathbf{0})$$

where  $\tilde{\mathbf{X}}$  is a  $N \times K$  matrix then

$$\tilde{\Sigma}_x = \begin{pmatrix} \hat{\Sigma}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$

and we can also augment  $\hat{\mathbf{B}}$  to a  $K \times K$  nonsingular matrix, say

$$\tilde{\mathbf{H}} = \begin{pmatrix} \hat{\mathbf{B}} \\ \mathbf{0} \end{pmatrix}$$

that still satisfies  $y = \tilde{x} \tilde{\mathbf{B}}$  and  $\hat{\Sigma}_y = \tilde{\mathbf{B}}^T \tilde{\Sigma}_x \tilde{\mathbf{B}}$ , then

$$\hat{\Sigma}_y^\dagger = \tilde{\mathbf{B}}^{-1} \begin{pmatrix} \hat{\Sigma}_x^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} (\tilde{\mathbf{B}}^T)^{-1},$$

then we have

$$\tilde{\mathbf{B}} \hat{\Sigma}_y^\dagger \tilde{\mathbf{B}}^T = \tilde{\Sigma}_x^\dagger,$$

then we conclude the proof.

## 7 ESL 4.5

Assume that we classify  $x_i$  to class 0 if  $x_i < x_0$ , the log-likelihood function is

$$l(\beta_0, \beta_1) = \sum_{i=1}^N [(\beta_0 + \beta_1 x_i) 1_{\{x_i > x_0\}} - \log(1 + e^{\beta_0 + \beta_1 x_i})],$$

then the MLE of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  satisfy

$$\frac{\partial l}{\partial \hat{\beta}_0} = \sum_{i=1}^N 1_{\{x_i > x_0\}} - \sum_{i=1}^N \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} = 0$$

and

$$\frac{\partial l}{\partial \hat{\beta}_1} = \sum_{i=1}^N x_i 1_{\{x_i > x_0\}} - \sum_{i=1}^N x_i \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} = 0,$$

moreover if consider that  $x_0$  separating two classes means that  $p(x_0; \beta_0, \beta_1) = 1/2$  then we have  $\beta_0 + x_0 \beta_1 = 0$ , then we substitute this into the one order equation and solve  $\beta$ .

Actually we can see the problem in another way, we want to maximize the probability of generating class 1 when  $x > x_0$ , note that this pobability is monotone decreasing for  $\beta_1$  with respect to all  $x_i > x_0$ , thus the value of MLE of  $\beta_0, \beta_1$  can be reached when  $\beta \rightarrow -\infty$ . In other words, the bigger  $|\beta_1|$  is the better when  $\beta_0 + x_0 \beta_1 = 0$ .

### 7.1 a

Assume that two classes can be separated by a hyperplane  $H_0$ , then for all  $x \in H_0$  there should be  $\beta_0 + x^T \beta_1 = 0$  hence  $H_0$  is exactly  $\beta_0 + x^T \beta_1 = 0$ , then we want to maximize the probability of generating class 1 when  $\beta_0 + x^T \beta_1 < 0$  which is

$$p(x; \beta_0, \beta_1) = \frac{1}{e^{\beta_0 + x^T \beta_1}},$$

this is monotone increasing for  $\|\beta\|$  if it maintains that  $\beta_0 + x^T \beta_1 < 0$ , or in other words  $\beta_0$  and  $\beta_1$  are mutiplied by the same constant. Hence the situation is similar with that in  $\mathbb{R}$ .

### 7.2 b

Assume that  $M$  classes are separated by the points  $-\infty = x_0 < x_1 < \dots < x_{M-1} < x_M = +\infty$  and if  $x_{m-1} < x < x_m$  we classify  $x$  to class  $m - 1$ . Similarly we should have

$$\beta_{m,0} + x_m \beta_{m,1} = \beta_{m+1,0} + x_m \beta_{m+1,1}, \quad m = 1, \dots, M - 2, \quad (7.2.1)$$

and

$$\beta_{M-1,0} + x_{M-1} \beta_{M-1,1} = 0, \quad (7.2.2)$$

also we want to maximize the probability of generating class  $M$  when  $x > x_{M-1}$  note that this pobability is monotone decreasing for  $\beta_{M,1}$  with respect to all  $x_i > x_{M-1}$ , then from equations (7.2.1) and euqation (7.2.2) we can obtain all  $\beta$ . Note that if we multiply all  $\beta$  by a constant bigger than 1 then the probability to generate all the sample points will increase, hence the same as the situation with 2 classes, the bigger all  $\beta$  the better when equations (7.2.1) and (7.2.2) hold.

## 8 ESL 4.6

### 8.1 a

By assumption if there is separability there exists a  $\beta$  such that  $\beta^T x_i^* > 0$  if  $y_i = 1$  and  $\beta^T x_i^* < 0$  if  $y_i = -1$  which can be written as  $y_i \beta^T x_i^* > 0$ , thus  $y_i \beta^T z_i > 0$ , hence there is a certain  $t$  such that  $y_i \beta^T z_i \geq t$  since the number of samples is finite, then let

$$\beta_{\text{sep}} = \frac{1}{m} \beta$$

that we conclude the proof.

### 8.2 b

We can calculate that

$$\|\beta_{\text{new}} - \beta_{\text{sep}}\|^2 = \|\beta_{\text{old}} - \beta_{\text{sep}} + y_i z_i\|^2 = \|\beta_{\text{old}} - \beta_{\text{sep}}\|^2 + y_i^2 \|z_i\|^2 + 2y_i(\beta_{\text{old}} - \beta_{\text{sep}})^T z_i,$$

hence we only need to verify that

$$y_i(\beta_{\text{old}} - \beta_{\text{sep}})^T z_i \leq -1$$

since  $y_i^2 \|z_i\|^2 = 1$ . This is quite trivial because  $y_i \beta_{\text{sep}}^T z_i = 1$  and  $y_i$  is misclassified thus  $y_i \beta_{\text{old}}^T z_i \leq 0$ . Thereby we conclude our proof.