

## Theoretical Questions

4.2 - a) The LDA decision rule is defined as

$$G(x) = \operatorname{argmax}_k \delta_k(x)$$

$$\delta(x) = x^\top \Sigma^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^\top \Sigma^{-1} \hat{\mu}_k + \log \pi_k \quad (4.10 \text{ in ESL})$$

Therefore  $G(x) = 2$  iff.  $\operatorname{argmax}_k (\delta_1, \delta_2) = 2$

$$\Leftrightarrow \delta_2 > \delta_1$$

$$\Leftrightarrow x^\top \Sigma^{-1} (\hat{\mu}_2 - \hat{\mu}_1) - \frac{1}{2} (\hat{\mu}_2 - \hat{\mu}_1)^\top \Sigma^{-1} (\hat{\mu}_2 - \hat{\mu}_1) + \log \left( \frac{N_2/N}{N_1/N} \right) > 0$$

$$\Leftrightarrow x^\top \Sigma^{-1} (\hat{\mu}_2 - \hat{\mu}_1) > \frac{1}{2} (\mu_2 - \mu_1)^\top \Sigma^{-1} (\hat{\mu}_2 - \hat{\mu}_1) + \log \left( \frac{N_2}{N_1} \right)$$

b)  $\hat{\beta} = (x^\top x)^{-1} x^\top Y$

$$X^\top X = \begin{pmatrix} 1_{N_1}^\top & 1_{N_2}^\top \\ X_1^\top & X_2^\top \end{pmatrix} \begin{pmatrix} 1_{N_1}^\top & X_1^\top \\ 1_{N_2}^\top & X_2^\top \end{pmatrix} = \begin{pmatrix} N_1 + N_2 & N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2 \\ N_1 \hat{\mu}_1^\top + N_2 \hat{\mu}_2^\top & X_1^\top X_1 + X_2^\top X_2 \end{pmatrix}$$

$$\sum_{k=1}^K \sum_{i=1}^{n_k} (x_i - \hat{\mu}_k) (x_i - \hat{\mu}_k)^\top = N_1 \hat{\mu}_1 \hat{\mu}_1^\top + N_2 \hat{\mu}_2 \hat{\mu}_2^\top$$

c)  $(\hat{\mu}_2 - \hat{\mu}_1)^\top$  is a  $1 \times p+1$  vector and  $\beta$  is a  $(p+1) \times 1$  vector, therefore  $(\hat{\mu}_2 - \hat{\mu}_1)^\top \beta$  is a scalar. Therefore  $\sum_B \beta$  is proportional to  $(\hat{\mu}_2 - \hat{\mu}_1)$ .

d)

4.3 - LDA using  $\hat{Y}$  is identical to LDA in the original space

iff.  $G_k(x) = G_k(y) \quad \forall k \in K$

let  $k$  be some class in the set  $K$ , and let

$$\delta_k(x) = x^T \Sigma^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \Sigma^{-1} \hat{\mu}_k$$

be the discriminant function for LDA with  $x$  as a predictor.

Want to prove that

$$\delta_k(x) \text{ is maximal} \Leftrightarrow \delta_k(y) \text{ is maximal} \quad (*)$$

$\Rightarrow$

let  $l \in K$  s.t.  $l \neq k$ . we have then

$$\delta_k(x) > \delta_l(x)$$

$$\Leftrightarrow x^T \Sigma^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \Sigma^{-1} \hat{\mu}_k + \log \pi_k > x^T \Sigma^{-1} \hat{\mu}_l - \frac{1}{2} \hat{\mu}_l^T \Sigma^{-1} \hat{\mu}_l + \log \pi_l$$

$$\Leftrightarrow \log\left(\frac{\pi_k}{\pi_l}\right) - \frac{1}{2} (\hat{\mu}_k + \hat{\mu}_l)^T \Sigma^{-1} (\hat{\mu}_k - \hat{\mu}_l) + x^T \Sigma^{-1} (\hat{\mu}_k - \hat{\mu}_l) > 0$$

Since  $\hat{\mu}_i = 1_{N_i}^T x_i / N_i$ , the sample means  $\hat{\mu}_i$  can be transformed by  $B$ , just like any other point in  $\mathbb{R}^P$ .

$$\Leftrightarrow -\frac{1}{2} (\hat{\mu}_k + \hat{\mu}_l)^T \Sigma^{-1} (\hat{\mu}_k - \hat{\mu}_l) + x^T \Sigma^{-1} (\hat{\mu}_k - \hat{\mu}_l) > 0$$

$$\Leftrightarrow x^T \Sigma^{-1} (\hat{\mu}_k - \hat{\mu}_l) > \frac{1}{2} (\hat{\mu}_k + \hat{\mu}_l)^T \Sigma^{-1} (\hat{\mu}_k - \hat{\mu}_l)$$

since multiplying both sides by  $(\hat{\mu}_k - \hat{\mu}_l)^T \Sigma^{-1} \Sigma$

$$\Leftrightarrow x^T > \frac{1}{2} (\hat{\mu}_k + \hat{\mu}_l)^T$$

$$\Leftrightarrow x^T \hat{B} > \frac{1}{2} (\hat{\mu}_k + \hat{\mu}_l)^T \hat{B}$$

$$\Leftrightarrow \hat{y} > \frac{1}{2} (\hat{\mu}_k' + \hat{\mu}_l')^T$$

where  $\hat{\mu}_i'$  is the average of the transformed  $\hat{y}_i$  sample:

$$\hat{\mu}_i' = \left[ 1_N^T x_i / N_i \right] \hat{B} = \frac{1}{N_i} 1_N^T Y$$

The other direction of (\*) can be obtained by applying  
a similar procedure as for  $\boxed{\Rightarrow}$ .