

# STOR 767 Spring 2019 Hw3

Due on 02/18/2019 in Class

Zhengan Fang

*Remark.* This homework focuses on linear classification methods.

## Instruction.

- Submission of handwritten solution for the **Theoretical Part** of this homework is allowed.
- Please use **RMarkdown** to create a formatted report for the **Computational Part** of this homework.
- Some of the problems are selected or modified from the textbook [1].

## Theoretical Part

1. (15 pt) (LDA and Least-Squared Rule, Textbook Ex. 4.2) Suppose we have training features  $\mathbf{X}_i \in \mathbb{R}^p$ , a two-class response with class sizes  $N_1, N_2$ ,  $N = N_1 + N_2$ , coded as

$$Y_i = \begin{cases} -N/N_1, & \text{the } i\text{-th sample belongs to class 1} \\ N/N_2, & \text{the } i\text{-th sample belongs to class 2} \end{cases}$$

- (a) Derive from scratch that the LDA rule classifies  $\mathbf{x} \in \mathbb{R}^p$  to class 2 if

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

and class 1 otherwise. Remember to declare  $\hat{\mu}_1, \hat{\mu}_2, \hat{\Sigma}$ .

*Hint.* Derive the Bayes rule under multivariate Gaussian with equal covariance assumption for class-conditional densities of  $\mathbf{X}_i$ 's, and plug in training sample estimates for unknown parameters.

The LDA rule classifies  $\mathbf{x}$  to class 2 if

$$\begin{aligned} & \log \frac{\Pr(G = 1 | X = x)}{\Pr(G = 2 | X = x)} < 0 \\ \iff & \log \frac{\pi_1}{\pi_2} - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + x^T \Sigma^{-1} (\mu_1 - \mu_2) < 0 \\ \iff & x^T \Sigma^{-1} (\mu_2 - \mu_1) > \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \log \pi_1 - \log \pi_2 \end{aligned}$$

Let  $\hat{\mu}_1, \hat{\mu}_2, \hat{\Sigma}, \frac{N_1}{N}, \frac{N_2}{N}$  be the estimation of  $\mu_1, \mu_2, \Sigma, \pi_1, \pi_2$  respectively, then we get the expression in the question.

(b) Consider minimization of the least squares criterion

$$\min_{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p} \sum_{i=1}^N (Y_i - \beta_0 - \beta^T \mathbf{X}_i)^2.$$

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

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

where  $\hat{\Sigma}_B := (\hat{\mu}_2 - \hat{\mu}_1)(\hat{\mu}_2 - \hat{\mu}_1)^T$ .

*Hint.* Solve  $\hat{\beta}_0(\beta)$  for fixed  $\beta$  first, then write down the normal equations for  $\hat{\beta}$ . Consider an ANOVA-type covariance matrix decomposition for  $\mathbf{X}_i$ 's:  $SS = SSW$  (within-class) +  $SSB$  (between-class). You might refer to Problem 2(b)(ii) for a more abstract development.

**Answer:**

Solve  $\hat{\beta}_0(\beta)$  for fixed  $\beta$ .

$$\hat{\beta}_0 = -\beta^T \frac{\sum_{i=1}^N \mathbf{X}_i}{N}$$

Write down the normal equations for  $\hat{\beta}$ .

$$\begin{aligned} \left( \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T \right) \hat{\beta} &= \sum_{i=1}^N (Y_i - \beta_0) \mathbf{X}_i \\ \iff \left( \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T \right) \hat{\beta} + \beta_0 \sum_{i=1}^N \mathbf{X}_i &= \sum_{i=1}^N Y_i \mathbf{X}_i \end{aligned}$$

Plug  $\hat{\beta}_0(\beta)$  in the above equation.

$$\begin{aligned} \left( \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T \right) \hat{\beta} - \left( \hat{\beta}^T \frac{\sum_{i=1}^N \mathbf{X}_i}{N} \right) \sum_{i=1}^N \mathbf{X}_i &= \sum_{i=1}^N Y_i \mathbf{X}_i \\ \iff \left( \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T - \frac{1}{N} \sum_{i=1}^N \mathbf{X}_i \sum_{i=1}^N \mathbf{X}_i^T \right) \hat{\beta} &= \sum_{i=1}^N Y_i \mathbf{X}_i. \end{aligned} \tag{a1}$$

By definition,

$$\begin{aligned} \hat{\mu}_1 &= \frac{1}{N_1} \sum_{g_i=1} \mathbf{X}_i \\ \hat{\mu}_2 &= \frac{1}{N_2} \sum_{g_i=2} \mathbf{X}_i \end{aligned}$$

$$\begin{aligned} \hat{\Sigma} &= \frac{1}{N-2} \left( \sum_{g_i=1} (\mathbf{X}_i - \hat{\mu}_1)(\mathbf{X}_i - \hat{\mu}_1)^T + \sum_{g_i=2} (\mathbf{X}_i - \hat{\mu}_2)(\mathbf{X}_i - \hat{\mu}_2)^T \right) \\ &= \frac{1}{N-2} \left( \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T - N_1 \hat{\mu}_1 \hat{\mu}_1^T - N_2 \hat{\mu}_2 \hat{\mu}_2^T \right) \\ \iff \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T &= (N-2)\hat{\Sigma} + N_1 \hat{\mu}_1 \hat{\mu}_1^T + N_2 \hat{\mu}_2 \hat{\mu}_2^T \end{aligned} \tag{a2}$$

$$\sum_{i=1}^N \mathbf{X}_i = N_1 \hat{\boldsymbol{\mu}}_1 + N_2 \hat{\boldsymbol{\mu}}_2 \quad (\text{a3})$$

$$\sum_{i=1}^N Y_i \mathbf{X}_i = \sum_{g_i=1} -N/N_1 \mathbf{X}_i + \sum_{g_i=2} N/N_2 \mathbf{X}_i = N(\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) \quad (\text{a4})$$

Plug (a2), (a3), (a4) into (a1), we get the expression in the question.

(c) Argue that

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

*i.e.* the least squares coefficient is identical to the LDA coefficient up to a scalar multiple.

Plug  $\hat{\boldsymbol{\Sigma}}_B = (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)(\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)^T$  into

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

we get

$$\begin{aligned} & \left( (N-2)\hat{\boldsymbol{\Sigma}} + \frac{N_1 N_2}{N} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)(\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)^T \right) \hat{\boldsymbol{\beta}} = N(\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) \\ \iff & (N-2)\hat{\boldsymbol{\Sigma}}\hat{\boldsymbol{\beta}} = \left( N - \frac{N_1 N_2}{N} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)^T \hat{\boldsymbol{\beta}} \right) (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) \\ \iff & \hat{\boldsymbol{\beta}} = \frac{1}{N-2} \left( N - \frac{N_1 N_2}{N} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)^T \hat{\boldsymbol{\beta}} \right) \hat{\boldsymbol{\Sigma}}^{-1}(\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) \\ \implies & \hat{\boldsymbol{\beta}} \propto \hat{\boldsymbol{\Sigma}}^{-1}(\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) \end{aligned}$$

(d) Argue that (1) holds invariantly for any (distinct) coding of  $Y_i$ 's.

For any distinct coding of  $Y_i$ ,

$$\sum_{i=1}^N Y_i \mathbf{X}_i = Y_1 \sum_{g_i=1} \mathbf{X}_i + Y_2 \sum_{g_i=2} \mathbf{X}_i = Y_1 N_1 \hat{\boldsymbol{\mu}}_1 + Y_2 N_2 \hat{\boldsymbol{\mu}}_2.$$

Then  $\hat{\boldsymbol{\beta}}$  satisfies

$$\left( (N-2)\hat{\boldsymbol{\Sigma}} + \frac{N_1 N_2}{N} \hat{\boldsymbol{\Sigma}}_B \right) \hat{\boldsymbol{\beta}} = Y_1 N_1 \hat{\boldsymbol{\mu}}_1 + Y_2 N_2 \hat{\boldsymbol{\mu}}_2$$

(e) Find the least-square prediction function  $\hat{f}_{\text{LS}}$ . Consider the least-square discriminant rule

$$\hat{d}_{\text{LS}}(\mathbf{x}) = \begin{cases} \text{class 2,} & \hat{f}_{\text{LS}}(\mathbf{x}) > 0 \\ \text{class 1,} & \text{otherwise} \end{cases}$$

Show that  $\hat{d}_{\text{LS}}$  is not the same as the LDA rules unless  $N_1 = N_2$ . If, alternatively, we use uniform class prior rather than plug-in estimates, then least-square discriminant rule is equivalent to the LDA.

Let

$$\hat{\beta} = \lambda \hat{\Sigma}^{-1}(\hat{\mu}_2 - \hat{\mu}_1)$$

where  $\lambda \in \mathbb{R}$  is a constant. Then,

$$\begin{aligned} \hat{\beta}_0 &= -\beta^T \frac{\sum_{i=1}^N \mathbf{X}_i}{N} \\ &= -\beta^T \frac{N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2}{N} \\ &= -\left(\lambda \hat{\Sigma}^{-1}(\hat{\mu}_2 - \hat{\mu}_1)\right)^T \frac{N_1 \hat{\mu}_1 + N_2 \hat{\mu}_2}{N} \\ &= -\frac{\lambda}{N} \left( (N_1 - N_2) \hat{\mu}_1 \hat{\Sigma}^{-1} \hat{\mu}_2^T - N_1 \hat{\mu}_1 \hat{\Sigma}^{-1} \hat{\mu}_1^T + N_2 \hat{\mu}_2 \hat{\Sigma}^{-1} \hat{\mu}_2^T \right) \end{aligned}$$

Then

$$\begin{aligned} \hat{f}_{\text{LS}} &= \hat{\beta}_0 + \hat{\beta}^T \mathbf{X}_i \\ &= -\frac{\lambda}{N} \left( (N_1 - N_2) \hat{\mu}_1 \hat{\Sigma}^{-1} \hat{\mu}_2^T - N_1 \hat{\mu}_1 \hat{\Sigma}^{-1} \hat{\mu}_1^T + N_2 \hat{\mu}_2 \hat{\Sigma}^{-1} \hat{\mu}_2^T \right) + \lambda \mathbf{X}_i^T \hat{\Sigma}^{-1}(\hat{\mu}_2 - \hat{\mu}_1) \end{aligned}$$

Therefore,  $\hat{d}_{\text{LS}}$  classifies  $\mathbf{X}_i$  to class 2 if

$$\begin{aligned} \hat{f}_{\text{LS}} &> 0 \\ \iff \mathbf{X}_i^T \hat{\Sigma}^{-1}(\hat{\mu}_2 - \hat{\mu}_1) &> \frac{1}{N} \left( (N_1 - N_2) \hat{\mu}_1 \hat{\Sigma}^{-1} \hat{\mu}_2^T - N_1 \hat{\mu}_1 \hat{\Sigma}^{-1} \hat{\mu}_1^T + N_2 \hat{\mu}_2 \hat{\Sigma}^{-1} \hat{\mu}_2^T \right) \end{aligned} \quad (\text{a5})$$

(a5) is not equivalent to the LDA rule (stated in question 1(a)) unless  $N_1 = N_2$ . Therefore,  $\hat{d}_{\text{LS}}$  is not the same as the LDA rules unless  $N_1 = N_2$ .

2. (20 pt) (Fisher's Discriminant Analysis/Reduced Rank LDA, Textbook Ex. 4.1, 4.8) Let  $\mathbf{X} = [X_{ij}]_{n \times p} \in \mathbb{R}^{n \times p}$  be a centered<sup>1</sup> data matrix with  $n$  observations (in rows) of  $p$ -dimensional features (in columns),  $\mathbf{Y} = [Y_{ik}]_{n \times K} \in \{0, 1\}^{n \times K}$  be an indicator matrix with 0-1 entries encoding memberships out of  $K (\leq n)$  classes for observations,  $n_1, n_2, \dots, n_K \geq 1$  be class sizes.

- (a) Derive the class-centroid matrix  $\mathbf{M} \in \mathbb{R}^{K \times p}$  in terms of  $\mathbf{X}$  and  $\mathbf{Y}$ .

$$\mathbf{Y}^T \mathbf{Y} = \text{diag}(n_1, n_2, \dots, n_K)$$

$$\mathbf{M} = (\text{diag}(n_1, n_2, \dots, n_K))^{-1} \mathbf{Y}^T \mathbf{X} = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X}$$

- (b) Denote rows in  $\mathbf{X}$  and  $\mathbf{M}$  as  $\{\mathbf{X}_{i\cdot}\}_{i=1}^n$  and  $\{\hat{\mu}_k\}_{k=1}^K$  respectively. Define

$$\mathbf{S} := \sum_{i=1}^n \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T, \quad \mathbf{S}_B := \sum_{k=1}^K n_k \hat{\mu}_k \hat{\mu}_k^T, \quad \mathbf{S}_W := \sum_{k=1}^K \sum_{Y_i=k} (\mathbf{X}_{i\cdot} - \hat{\mu}_k)(\mathbf{X}_{i\cdot} - \hat{\mu}_k)^T.$$

---

<sup>1</sup>That is,  $\sum_{i=1}^n X_{ij} = 0$  for all  $1 \leq j \leq p$ .

- (i) Write  $\mathbf{S}, \mathbf{S}_B, \mathbf{S}_W$  in compact form in terms of  $\mathbf{X}$  and  $\mathbf{Y}$ .

$$\mathbf{S} = \mathbf{X}^T \mathbf{X}$$

$$\mathbf{S}_B = \mathbf{M}^T \text{diag}(n_1, \dots, n_K) \mathbf{M} = \mathbf{M}^T \mathbf{Y}^T \mathbf{Y} \mathbf{M} = \mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X}$$

$$\begin{aligned} \mathbf{S}_W &= (\mathbf{X} - \mathbf{Y} \mathbf{M})^T (\mathbf{X} - \mathbf{Y} \mathbf{M}) \\ &= \mathbf{X}^T \mathbf{X} - \mathbf{M}^T \mathbf{Y}^T \mathbf{X} - \mathbf{X}^T \mathbf{Y} \mathbf{M} + \mathbf{M}^T \mathbf{Y}^T \mathbf{Y} \mathbf{M} \\ &= \mathbf{X}^T \mathbf{X} - 2\mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} + \mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \\ &= \mathbf{X}^T \mathbf{X} - \mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \end{aligned}$$

- (ii) Prove an ANOVA-type covariance matrix decomposition

$$\mathbf{S} = \mathbf{S}_B + \mathbf{S}_W.$$

Plug in the results from (i).

- (c) Suppose  $\mathbf{S}_W > 0$ . Show how to solve the generalized eigenvalue problem

$$\begin{aligned} \max_{\boldsymbol{\phi}} \quad & \boldsymbol{\phi}^T \mathbf{S}_B \boldsymbol{\phi} \\ \text{s.t.} \quad & \boldsymbol{\phi}^T \mathbf{S}_W \boldsymbol{\phi} = 1 \end{aligned} \tag{2}$$

by transforming it to a standard eigenvalue problem.

There is an orthogonal matrix  $\mathbf{P}$  and a diagonal matrix  $\mathbf{D}$  such that  $\mathbf{P}^T \mathbf{S}_W \mathbf{P} = \mathbf{D}$ . Let  $\mathbf{S}_W^{\frac{1}{2}} = \mathbf{P} \mathbf{D}^{\frac{1}{2}} \mathbf{P}^T$ ,  $\mathbf{S}_W^{-\frac{1}{2}} = \mathbf{P} \mathbf{D}^{-\frac{1}{2}} \mathbf{P}^T$ ,  $\boldsymbol{\theta} = \mathbf{S}_W^{\frac{1}{2}} \boldsymbol{\phi}$ . Transform (2) to

$$\begin{aligned} \max_{\boldsymbol{\theta}} \quad & \boldsymbol{\theta}^T \mathbf{S}_W^{-\frac{1}{2}} \mathbf{S}_B \mathbf{S}_W^{-\frac{1}{2}} \boldsymbol{\theta} \\ \text{s.t.} \quad & \boldsymbol{\theta}^T \boldsymbol{\theta} = 1 \end{aligned}$$

- (d) Suppose  $\mathbf{S}_W > 0$ . Let  $\{\hat{\boldsymbol{\phi}}_j\}_{j=1}^p$  be Fisher's discriminant coordinates, *i.e.* the solutions to (2) with each one orthogonal in  $\mathbf{S}_W$ <sup>2</sup> to the predecessors. For  $1 \leq m \leq p$ , define the (Fisher's) discriminant rules based on nearest centroids projected on  $\{\hat{\boldsymbol{\phi}}_j\}_{j=1}^m$

$$d_m(\mathbf{x}) = \underset{1 \leq k \leq K}{\operatorname{argmin}} \sum_{j=1}^m \left( \hat{\boldsymbol{\phi}}_j^T \mathbf{x} - \hat{\boldsymbol{\phi}}_j^T \hat{\boldsymbol{\mu}}_k \right)^2.$$

Now let  $L := \dim(\operatorname{span}\{\hat{\boldsymbol{\mu}}_k\}_{k=1}^K) = \operatorname{rank}(\mathbf{S}_B) \leq \min\{K-1, p\}$ . Take the following steps to show the equivalence between LDA based on multivariate-Gaussian Bayes rule introduced in class and Fisher's discrimination.

- (i) Using the fact that  $\hat{\boldsymbol{\phi}}_j^T \mathbf{S}_W \hat{\boldsymbol{\phi}}_{j'} = \mathbb{1}(j = j')$  ( $1 \leq j, j' \leq p$ ), show that

$$\mathbf{S}_W^{-1} = \sum_{j=1}^p \hat{\boldsymbol{\phi}}_j \hat{\boldsymbol{\phi}}_j^T.$$

---

<sup>2</sup>We say  $\mathbf{x}$  is orthogonal to  $\mathbf{y}$  in  $\mathbf{A}$  if  $\mathbf{x}^T \mathbf{A} \mathbf{y} = 0$ .

Let

$$\mathbf{A} = \begin{bmatrix} \hat{\phi}_1^T \\ \vdots \\ \hat{\phi}_p^T \end{bmatrix}$$

Then,

$$\begin{aligned} \hat{\phi}_j^T \mathbf{S}_W \hat{\phi}_{j'} &= \mathbb{1}(j = j') \\ \implies \mathbf{A} \mathbf{S}_W \mathbf{A}^T &= \mathbf{I}_p \\ \implies \mathbf{A} \mathbf{S}_W &= (\mathbf{A}^T)^{-1} \\ \implies \mathbf{A}^T \mathbf{A} \mathbf{S}_W &= \mathbf{I}_p \\ \implies \mathbf{S}_W^{-1} &= \mathbf{A}^T \mathbf{A} = \sum_{j=1}^p \hat{\phi}_j \hat{\phi}_j^T \end{aligned}$$

(ii) Show that

$$\hat{\phi}_j^T \mathbf{S}_B \hat{\phi}_j = 0. \quad (L+1 \leq j \leq p)$$

Because  $\mathbf{S}_W^{-\frac{1}{2}}$  has full rank,  $\mathbf{S}_W^{-\frac{1}{2}} \mathbf{S}_B \mathbf{S}_W^{-\frac{1}{2}}$  has rank of  $L$ . Therefore,  $\mathbf{S}_W^{-\frac{1}{2}} \mathbf{S}_B \mathbf{S}_W^{-\frac{1}{2}}$  has  $p-L$  zero eigenvalues. Therefore,  $\hat{\phi}_j^T \mathbf{S}_B \hat{\phi}_j = 0. \quad (L+1 \leq j \leq p)$ .

(iii) Show the equivalence between (ii) and

$$\hat{\phi}_j^T \hat{\mu}_k = 0. \quad (L+1 \leq j \leq p, 1 \leq k \leq K)$$

$$\begin{aligned} \hat{\phi}_j^T \mathbf{S}_B \hat{\phi}_j &= 0 \\ \iff \sum_{k=1}^K n_k \hat{\phi}_j^T \hat{\mu}_k \hat{\mu}_k^T \hat{\phi}_j &= 0 \\ \iff \sum_{k=1}^K n_k (\hat{\phi}_j^T \hat{\mu}_k)^2 &= 0 \\ \iff \hat{\phi}_j^T \hat{\mu}_k &= 0 \quad (1 \leq k \leq K) \end{aligned}$$

(iv) Show that  $d_L, d_{L+1}, \dots, d_p$  are equivalent to the multivariate-Gaussian-based LDA with uniform prior.

Let  $d(\mathbf{x})$  be the estimation from multivariate-Gaussian-based LDA with uniform prior.

$$\begin{aligned}
d(\mathbf{x}) &= \operatorname{argmin}_{1 \leq k \leq K} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k) \\
&= \operatorname{argmin}_{1 \leq k \leq K} \frac{1}{n - K} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k)^T \mathbf{S}_W^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k) \\
&= \operatorname{argmin}_{1 \leq k \leq K} \sum_{j=1}^p (\mathbf{x} - \hat{\boldsymbol{\mu}}_k)^T \hat{\boldsymbol{\phi}}_j \hat{\boldsymbol{\phi}}_j^T (\mathbf{x} - \hat{\boldsymbol{\mu}}_k) \\
&= \operatorname{argmin}_{1 \leq k \leq K} \sum_{j=1}^p (\hat{\boldsymbol{\phi}}_j^T \mathbf{x} - \hat{\boldsymbol{\phi}}_j^T \hat{\boldsymbol{\mu}}_k)^2 \\
&= d_p(\mathbf{x}).
\end{aligned}$$

$$\begin{aligned}
\hat{\boldsymbol{\phi}}_j^T \hat{\boldsymbol{\mu}}_k &= 0 \quad (L+1 \leq j \leq p, 1 \leq k \leq K) \\
\implies d_L(\mathbf{x}) &= d_{L+1}(\mathbf{x}) = \dots = d_p(\mathbf{x}).
\end{aligned}$$

Therefore,  $d_L, d_{L+1}, \dots, d_p$  are equivalent to the multivariate-Gaussian-based LDA with uniform prior.

*Remark* (Reduced-Rank LDA). The scaling matrices for Fisher's discriminant rules are actually [low-rank approximations](#)<sup>3</sup> to that of multivariate-Gaussian LDA. In particular,  $m = p$  recovers the scaling matrix of multivariate-Gaussian LDA as in (i), and the  $m(\geq L)$ -rank approximations don't lose information in the sense of discriminating centroids in the reduced space as in (iii). The equivalence of (ii) and (iii) reveals the rationale of finding directions to maximize between-class variances.

3. (15 pt) (Logistic Regression Recession, Textbook Ex. 4.5) Consider a two-class Logistic regression problem. We encode labels  $Y_i \in \{\pm 1\}$ .

(a) Show that solving the maximum-likelihood estimate (MLE) for the intercept and slope parameters  $(\beta_0, \boldsymbol{\beta})$  is equivalent to the following empirical risk minimization (ERM) problem<sup>4</sup>

$$\min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{n} \sum_{i=1}^n \log \{1 + \exp[-Y_i(\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta})]\}.$$

**Punchline.** In standard generalized linear model (GLM) formulation<sup>5</sup> for the Logistic likelihood,  $Y_i$ 's are coded as 0-1, but now we are coding it into  $\pm 1$ .

Let

$$y_i = \begin{cases} 1 & Y_i = 1 \\ 0 & Y_i = -1 \end{cases}$$

<sup>3</sup>NOT under the usual spectral/Frobenius norm, but under a Mahalanobis-type Frobenius norm.

<sup>4</sup>We call  $\ell(u, v) := \log(1 + e^{-uv})$  as the **Logistic loss**.

<sup>5</sup>See Section 4.4.1 at page 114 in [2, McCullagh and Nelder (1989)].

The maximum-likelihood estimate is

$$\begin{aligned}
& \max_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^n \left\{ y_i \log \frac{\exp[\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta}]}{1 + \exp[\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta}]} + (1 - y_i) \log \frac{1}{1 + \exp[\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta}]} \right\} \\
& \iff \max_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^n \left\{ y_i \log \frac{1}{1 + \exp[-(\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta})]} + (1 - y_i) \log \frac{1}{1 + \exp[\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta}]} \right\} \\
& \iff \min_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^n \{ y_i \log \{1 + \exp[-(\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta})]\} + (1 - y_i) \log \{1 + \exp[\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta}]\} \} \\
& \iff \min_{(\beta_0, \boldsymbol{\beta})} \sum_{i=1}^n \log \{1 + \exp[-Y_i(\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta})]\}
\end{aligned}$$

- (b) Argue that if the training data can be perfectly separated by a hyperplane (Figure 4.16 at page 134 in Textbook), or mathematically, there exists a hyperplane  $H \subseteq \mathbb{R}^p$  with upper (*resp.* lower) halfspace  $H^+$  (*resp.*  $H^-$ )<sup>6</sup> such that

$$\{\mathbf{X}_i : Y_i = +1 \ (1 \leq i \leq n)\} \subseteq H^+, \quad \{\mathbf{X}_i : Y_i = -1 \ (1 \leq i \leq n)\} \subseteq H^-,$$

then the optimal solution to Logistic regression is characterized by some unbounded directions<sup>7</sup>. That is, there exists a direction  $(\alpha, \boldsymbol{\gamma}) \in \mathbb{R}^{p+1}$  along which  $(\beta_0, \boldsymbol{\beta}) = c \times (\alpha, \boldsymbol{\gamma})$  will give objectives increasing/decreasing to an objective supremum/infimum as  $c \rightarrow +\infty$ .

Let  $H = \alpha + \mathbf{X}^T \boldsymbol{\gamma}$ . Then,

$$\begin{aligned}
\alpha + \mathbf{X}_i^T \boldsymbol{\gamma} &> 0, \quad \text{if } Y_i = +1 \\
\alpha + \mathbf{X}_i^T \boldsymbol{\gamma} &< 0, \quad \text{if } Y_i = -1
\end{aligned}$$

Therefore,

$$Y_i(\alpha + \mathbf{X}_i^T \boldsymbol{\gamma}) > 0, \quad 1 \leq i \leq n$$

Then,

$$\sum_{i=1}^n \log \{1 + \exp[-Y_i(\beta_0 + \mathbf{X}_i^T \boldsymbol{\beta})]\} = \sum_{i=1}^n \log \{1 + \exp[-cY_i(\alpha + \mathbf{X}_i^T \boldsymbol{\gamma})]\} \rightarrow 0$$

as  $c \rightarrow +\infty$ .

- (c) Describe the recession/unboundedness phenomenon in (b) when the number of classes is greater than two.

Denote the number of classes as  $K$ . Denote the class of the  $i$ th sample as  $G_i$ .

If for each  $k$ ,  $1 \leq k \leq K-1$ , there exists a hyperplane  $H \subseteq \mathbb{R}^p$  with upper (*resp.* lower) halfspace  $H^+$  (*resp.*  $H^-$ ) such that

$$\{\mathbf{X}_i : G_i = k \ (1 \leq i \leq n)\} \subseteq H^+, \quad \{\mathbf{X}_i : G_i \neq k \ (1 \leq i \leq n)\} \subseteq H^-,$$

<sup>6</sup>If we write a hyperplane into  $H = \{\mathbf{x} \in \mathbb{R}^p : \beta_0 + \mathbf{x}^T \boldsymbol{\beta} = 0\}$ , then the halfspaces are  $H^{+(-)} = \{\mathbf{x} \in \mathbb{R}^p : \beta_0 + \mathbf{x}^T \boldsymbol{\beta} \geq (\leq) 0\}$ .

<sup>7</sup>In convex analysis, we say such directions compose of the **recession cone** of the objective function. More general discussion of recession and unboundedness can be found in Section 8 in [3, Rockafellar (1970)].



then the optimal solution to Logistic regression is characterized by some unbounded directions. That is, there exists directions  $(\alpha_k, \gamma_k) \in \mathbb{R}^{p+1}$  ( $1 \leq k \leq K-1$ ) along which  $(\beta_{k0}, \beta_k) = c \times (\alpha_k, \gamma_k)$  will give objectives increasing/decreasing to an objective supremum/infimum as  $c \rightarrow +\infty$ .

## Computational Part

### 1. (25 pt) South African Heart Disease Data

Perform LDA, QDA, and Logistic regression analysis on the *South African Heart Disease Data*<sup>8</sup> and compare their prediction performances. Find more information in Textbook Section 4.4.2. Remember to hold out a test set to perform prediction.

### 2. (25 pt) Zip Code Digits Data

Classify between 3's and 8's from the *Zip Code Digits* data<sup>9</sup>. Build up the LDA, QDA, and Logistic regression models on the training dataset and compare their prediction performances on the test dataset.

## References

- [1] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*. Springer-Verlag, <https://web.stanford.edu/~hastie/Papers/ESLII.pdf>, second edition, 2009. 1
- [2] P. McCullagh and J.A. Nelder. *Generalized linear models*. Chapman and Hall, second edition, 1989. 7
- [3] Ralph Tyrell Rockafellar. *Convex analysis*. Princeton university press, 1970. 8

---

<sup>8</sup>Available at <https://web.stanford.edu/~hastie/ElemStatLearn/datasets/SAheart.data>.

<sup>9</sup>Available in <https://web.stanford.edu/~hastie/ElemStatLearn/datasets/zip.train.gz> (training dataset) and <https://web.stanford.edu/~hastie/ElemStatLearn/datasets/zip.test.gz> (test dataset).