

## Homework 7 (Due 6/5)

**Question 1** Suppose the population mean and covariance matrix of  $\begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$  are

$$\vec{\mu} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 2 & \sqrt{3} \\ \sqrt{3} & 4 \end{bmatrix}.$$

- (a) Determine the first and second principal components  $Y_1$  and  $Y_2$ , and find their variances, respectively.
- (b) Determine the proportion of total variance due to the  $Y_1$ .
- (c) Compare the contributions of  $X_1$  and  $X_2$  to the determination of  $Y_1$  based on loadings and correlations, respectively.

**Question 2** Suppose that the random vector  $\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$  has the following population mean and covariance matrix:

$$\mu = \mathbf{0}, \quad \Sigma = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 1 \\ 0 & 1 & 3 \end{bmatrix}.$$

- (a) Determine the first two principal components  $Y_1$  and  $Y_2$  and the proportion of total variance due to them. For  $i = 1, 2$ , compare the contributions of  $X_1$ ,  $X_2$  and  $X_3$  to the determination of  $Y_i$  based on loadings and correlations, respectively.
- (b) Let  $Z_1, Z_2, Z_3$  be the standardized variables of  $X_1, X_2, X_3$ , respectively. Find the first two principal components  $W_1$  and  $W_2$  of  $(Z_1, Z_2, Z_3)$ . For  $i = 1, 2$ , compare the contributions of  $Z_1$ ,  $Z_2$  and  $Z_3$  to the determination of  $W_i$  based on loadings and correlations, respectively.

**Question 3** Suppose a sample  $\vec{x}_1, \dots, \vec{x}_n$  has the sample mean  $\bar{\vec{x}}$  and sample covariance  $\mathbf{S}$ . Let  $\hat{y}_i$  be the  $i$ -th sample principal component, where  $i = 1, \dots, p$ . Then we transform the data matrix

$$\mathbf{X} = \begin{bmatrix} \vec{x}_1^\top \\ \vec{x}_2^\top \\ \vdots \\ \vec{x}_n^\top \end{bmatrix}$$

into the data matrix of the first  $r (< p)$  sample principal components

$$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{y}_{11} & \hat{y}_{12} & \cdots & \hat{y}_{1r} \\ \hat{y}_{21} & \hat{y}_{22} & \cdots & \hat{y}_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{y}_{n1} & \hat{y}_{n2} & \cdots & \hat{y}_{nr} \end{bmatrix}.$$

Please find the relationship between  $\mathbf{X}$  and  $\hat{\mathbf{Y}}$  with the spectral decomposition of  $\mathbf{S}$ .

**Question 4** You are given a sample  $\vec{x}_1, \dots, \vec{x}_6$  from a 2-dimension population. Moreover, the sample principal components are

$$\hat{y}_{i1} = \frac{\sqrt{2}}{2}x_{i1} + \frac{\sqrt{2}}{2}x_{i2}, \quad \hat{y}_{i2} = \frac{\sqrt{2}}{2}x_{i1} - \frac{\sqrt{2}}{2}x_{i2}, \quad i = 1 \dots, 6.$$

Assume the sample variances of  $(\hat{y}_{i1})_{i=1}^6$  and  $(\hat{y}_{i2})_{i=1}^6$  are 3 and 2, respectively. Find the sample covariance matrix of the original data matrix  $\mathbf{X}$ .

**Question 5** Consider a sample  $\vec{x}_{11}, \dots, \vec{x}_{1n_1}$  of size  $n_1 = 10$  from population 1 (corresponding to class  $\pi_1$ ) and a sample  $\vec{x}_{21}, \dots, \vec{x}_{2n_2}$  of size  $n_2 = 10$  from population 2 (corresponding to class  $\pi_2$ ). The summary statistics for these two samples are

$$\begin{aligned} \bar{\vec{x}}_1 &= \begin{bmatrix} 6 \\ 0 \end{bmatrix}, & \mathbf{S}_1 &= \begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}, \\ \bar{\vec{x}}_2 &= \begin{bmatrix} 0 \\ 3 \end{bmatrix}, & \mathbf{S}_2 &= \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix}. \end{aligned}$$

For some  $\vec{x}_0 = \begin{bmatrix} x_{01} \\ x_{02} \end{bmatrix}$ , derive the following classifiers:

1. Classifier 1: Fisher's rule based on  $\vec{x}_0$ , and give the  $D^2$  distance denoted as  $D_1^2$ ;
2. Classifier 2: Fisher's rule based on  $x_{01}$ , and give the  $D^2$  distance denoted as  $D_2^2$ ;
3. Classifier 3: Fisher's rule based on  $x_{02}$ , and give the  $D^2$  distance denoted as  $D_3^2$ ;
4. Classifier 4: Fisher's rule based on  $\vec{x}_0^\top (\bar{\vec{x}}_1 - \bar{\vec{x}}_2)$ , and give the  $D^2$  distance denoted as  $D_4^2$ ;
5. Classifier 5: Fisher's rule based on  $\vec{x}_0^\top \mathbf{S}_{pooled}^{-1} (\bar{\vec{x}}_1 - \bar{\vec{x}}_2)$ , and give the  $D^2$  distance denoted as  $D_5^2$ ;
6. Compare the above  $D^2$  distances.

**Question 6** Suppose we have  $n_1$   $p$ -variate observations from  $\pi_1$  and  $n_2$   $p$ -variate observations from  $\pi_2$ . The respective data matrices are

$$\mathbf{X}_1 = \begin{bmatrix} \vec{x}_{11}^\top \\ \vec{x}_{12}^\top \\ \vdots \\ \vec{x}_{1n_1}^\top \end{bmatrix}, \quad \mathbf{X}_2 = \begin{bmatrix} \vec{x}_{21}^\top \\ \vec{x}_{22}^\top \\ \vdots \\ \vec{x}_{2n_2}^\top \end{bmatrix}.$$

Suppose  $\bar{\vec{x}}$  is the overall sample mean of these two samples. Consider the data matrices

$$\mathbf{Z}_1 = \begin{bmatrix} \vec{z}_{11}^\top \\ \vec{z}_{12}^\top \\ \vdots \\ \vec{z}_{1n_1}^\top \end{bmatrix}, \quad \mathbf{Z}_2 = \begin{bmatrix} \vec{z}_{21}^\top \\ \vec{z}_{22}^\top \\ \vdots \\ \vec{z}_{2n_2}^\top \end{bmatrix},$$

where  $\vec{z}_{lj} = \vec{x}_{lj} - \bar{\vec{x}}$ ,  $l = 1, 2$ ,  $j = 1, \dots, n_l$ . Show that by Fisher's linear discriminant,  $\vec{x}_0$  is allocated to the first population based on  $\mathbf{X}_1$  and  $\mathbf{X}_2$  if and only if  $\vec{x}_0 - \bar{\vec{x}}$  is allocated to the first population based on  $\mathbf{Z}_1$  and  $\mathbf{Z}_2$ .

**Question 7** Consider three independent samples from three classes:

$\pi_1$ : distribution  $\mathcal{N}_p(\vec{\mu}_1, \mathbf{\Sigma})$ , sample size  $n_1$ , sample mean  $\vec{x}_1$ , sample covariance  $\mathbf{S}_1$ ;

$\pi_2$ : distribution  $\mathcal{N}_p(\vec{\mu}_2, \mathbf{\Sigma})$ , sample size  $n_2$ , sample mean  $\vec{x}_2$ , sample covariance  $\mathbf{S}_2$ ;

$\pi_3$ : distribution  $\mathcal{N}_p(\vec{\mu}_3, \mathbf{\Sigma})$ , sample size  $n_3$ , sample mean  $\vec{x}_3$ , sample covariance  $\mathbf{S}_3$ .

Assume that these three populations have equal prior probabilities. For a new observation  $\vec{x}_0$ , we aim to classify it to one of the three classes with pairwise linear discriminant analyses based on Fisher's rule. Given the three population covariance matrices are assumed to be the same, in all linear discriminant analyses we use

$$\mathbf{S}_{pooled} = \frac{1}{n_1 + n_2 + n_3 - 3} ((n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2 + (n_3 - 1)\mathbf{S}_3).$$

Suppose that in the comparison between  $\pi_1$  and  $\pi_2$ ,  $\vec{x}_0$  is allocated to  $\pi_2$ , while in the comparison between  $\pi_2$  and  $\pi_3$ ,  $\vec{x}_0$  is allocated to  $\pi_3$ . Show that in the comparison between  $\pi_1$  and  $\pi_3$ ,  $\vec{x}_0$  is allocated to  $\pi_3$ .