



Departamento de Matemática

## Multivariate Analysis

Mater in Eng. and Data Science & Master in Mathematics and Applications

2<sup>nd</sup> Test

Duration: 1.5 hours

1<sup>st</sup> Semester – 2020/2021

17/12/2020 – 18:30

Please justify conveniently your answers

---

### Group I

**9.0 points**

1. Let  $\mathbf{x} = (x_1, \dots, x_p)^t$  and  $\mathbf{y} = (y_1, \dots, y_p)^t$  represents the  $p$  continuous measurements characterizing two different objects, where for  $\mathbf{x}, \mathbf{y} \in (\mathbb{R}^+)^p$ ,

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\frac{1}{p} \sum_{i=1}^p \left( \frac{x_i - y_i}{x_i + y_i} \right)^2} \quad \text{and} \quad d(\mathbf{0}, \mathbf{0}) = 0.$$

Prove that the previous coefficient is a dissimilarity measure.

2. Show that the distance  $d_{k(ij)}$  (between cluster  $k$  and the cluster formed by merging cluster  $i$  and  $j$ ) used by average-linkage verifies:

$$d_{k(ij)} = \frac{n_i}{n_i + n_j} d_{ki} + \frac{n_j}{n_i + n_j} d_{kj},$$

where  $d_{ij}$  is the dissimilarity between  $i$ -th and  $j$ -th cluster and  $n_i$  is the number of objects belonging to the  $i$ -th cluster.

3. The pairwise dissimilarities between four objects are as follows:

$$\mathcal{D} = \begin{bmatrix} 0 & 1 & 11 & 5 \\ 1 & 0 & 2 & 3 \\ 11 & 2 & 0 & 4 \\ 5 & 3 & 4 & 0 \end{bmatrix}.$$

Use average-linkage cluster analysis on the dissimilarity matrix above, and draw the associated dendrogram.

**Suggestion:** Use the result stated in Question 2.

---

### Group II

**11.0 points**

1. An observation  $x$  comes from one of two populations with prior probabilities  $P(Y = 0) = P(Y = 1)$  and probability mass functions:

$$P(X = x|Y = j) = \begin{cases} \frac{\lambda_j^x e^{-\lambda_j}}{x!}, & x = 0, 1, 2, \dots \\ 0, & \text{otherwise,} \end{cases}$$

with  $\lambda_1 > \lambda_0 > 0$ ,  $j = 0, 1$ , and  $X|Y = j \sim \text{Poisson}(\lambda_j)$ .

- (a) Obtain the classification rule that minimizes the total probability of misclassification. (4.0)

- (b) Assuming  $\lambda_0 = 1$ , and  $\lambda_1 = 4$ , obtain:

i. The recall of class 1. (2.5)

ii. The precision of class 1. (2.5)

iii. The  $F_1$  measure of class 1. (1.0)

- (c) Comment on the accuracy of the classifier. (1.0)



Departamento de Matemática

## Multivariate Analysis

Mater in Eng. and Data Science & Master in Mathematics and Applications

2<sup>nd</sup> Test

1<sup>st</sup> Semester – 2020/2021

Duration: 1.5 hours

17/12/2020 – 18:30

Please justify conveniently your answers

---

### Group I

9.0 points

1. Let  $\underline{x} = (x_1, \dots, x_p)^t$  and  $\underline{y} = (y_1, \dots, y_p)^t$  represents the  $p$  continuous measurements characterizing two different objects, where for  $\underline{x}, \underline{y} \in (\mathbb{R}^+)^p$ , (3.0)

$$d(\underline{x}, \underline{y}) = \sqrt{\frac{1}{p} \sum_{i=1}^p \left( \frac{x_i - y_i}{x_i + y_i} \right)^2} \quad \text{and} \quad d(\mathbf{0}, \mathbf{0}) = 0.$$

Prove that the previous coefficient is a dissimilarity measure.

We have to prove:

1.  $d(\underline{x}, \underline{x}) = 0$
2.  $d(\underline{x}, \underline{y}) \geq 0 \quad \forall \underline{x}, \underline{y} \in \mathbb{R}^p$
3.  $d(\underline{x}, \underline{y}) = d(\underline{y}, \underline{x})$

Proof:

$$1. \text{ If } \underline{x} \neq \underline{0} \Rightarrow d(\underline{x}, \underline{x}) = \sqrt{\frac{1}{p} \sum_{i=1}^p \left( \frac{x_i - x_i}{x_i + x_i} \right)^2} = 0$$

If  $\underline{x} = \underline{0} \Rightarrow d(\underline{0}, \underline{0}) = 0$  so (1) is proved.

$$2. \quad d(\underline{x}, \underline{y}) = \sqrt{\frac{1}{p} \sum_{i=1}^p \left( \frac{x_i - y_i}{x_i + y_i} \right)^2} \geq 0$$

so  $d(\underline{x}, \underline{y}) \geq 0$

And finally,

$$\begin{aligned} (3) \quad d(\underline{y}, \underline{x}) &= \sqrt{\frac{1}{p} \sum_{i=1}^p \left( \frac{y_i - x_i}{y_i + x_i} \right)^2} = \\ &= \sqrt{\frac{1}{p} \sum_{i=1}^p \left( \frac{(y_i - x_i)^2}{(y_i + x_i)^2} \right)} = d(\underline{x}, \underline{y}) \end{aligned}$$

QED

2. Show that the distance  $d_{k(ij)}$  (between cluster  $k$  and the cluster formed by merging cluster  $i$  and  $j$ ) used by average linkage verifies:

$$d_{k(ij)} = \frac{n_i}{n_i + n_j} d_{ki} + \frac{n_j}{n_i + n_j} d_{kj},$$

where  $d_{ij}$  is the dissimilarity between  $i$ -th and  $j$ -th cluster and  $n_i$  is the number of objects belonging to the  $i$ -th cluster.

$$\begin{aligned}
 d_{k(ij)} &= \frac{1}{(n_k)(n_i+n_j)} \sum_{u=1}^{n_k} \sum_{v=1}^{n_i+n_j} d_{uv} = \\
 &= \frac{1}{n_k(n_i+n_j)} \sum_{u=1}^{n_k} \left[ \sum_{v=1}^{n_i} d_{uv} + \sum_{v=1}^{n_j} d_{uv} \right] \\
 &= \frac{n_i}{n_i+n_j} \sum_{u=1}^{n_k} \sum_{v=1}^{n_i} \frac{d_{uv}}{n_k n_i} + \frac{n_j}{n_i+n_j} \sum_{u=1}^{n_k} \sum_{v=1}^{n_j} d_{uv} \\
 &= \frac{n_i}{n_i+n_j} d_{ki} + \frac{n_j}{n_i+n_j} d_{kj} \quad \text{QED}
 \end{aligned}$$

3. The pairwise dissimilarities between four objects are as follows:

$$D = \begin{bmatrix} 0 & 1 & 11 & 5 \\ 1 & 0 & 2 & 3 \\ 11 & 2 & 0 & 4 \\ 5 & 3 & 4 & 0 \end{bmatrix}.$$

Use average-linkage cluster analysis on the dissimilarity matrix above, and draw the associated dendrogram. (4.0)

**Suggestion:** Use the result state in Question 2.

1st step: when  $d_{ij} = d_{12} = 1$

$$D = \begin{bmatrix} 0 & 1 & 11 & 5 \\ 1 & 0 & 2 & 3 \\ 11 & 2 & 0 & 4 \\ 5 & 3 & 4 & 0 \end{bmatrix}. \quad d(3, \{1, 2\}) = \frac{1}{1+1} d(1, 3) + d(2, 3) \frac{1}{1+1} = \frac{11+2}{2} = 6.5$$

$$d(4, \{1, 2\}) = \frac{1}{1+1} d(4, 1) + d(4, 2) \frac{1}{2} = \frac{5+3}{2} = 4$$

thus, the new dissimilarity matrix is:

$$\begin{array}{c} \{1, 2\} \quad 3 \quad 4 \\ \begin{array}{c} 0 \\ 6.5 \\ 4 \end{array} \end{array}$$

2nd step:

when  $d_{ij} = d(\{1, 2\}, 4) = d(3, 4) = 4$

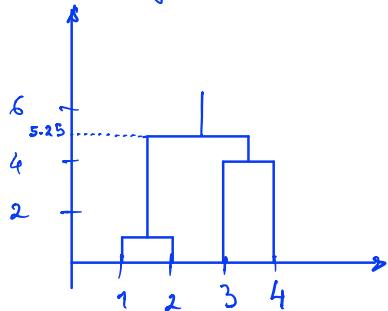
$$d(\{1, 2\}, \{3, 4\}) = d(\{1, 2\}, 3) \frac{1}{1+1} +$$

$$+ d(\{1, 2\}, 4) = \frac{1}{1+1} = \frac{6.5 + 4}{2} = \frac{10.5}{2} = 5.25$$

$$D = \begin{bmatrix} \{1, 2\} & 0 \\ \{3, 4\} & 5.25 \end{bmatrix}$$

$$3^{\text{rd}} \text{ step: } \text{new } d_{ij} = \boxed{d(\{1, 2\}, \{3, 4\}) = 5.25}$$

Dendrogram:



## 2<sup>nd</sup> Possible Solution

1<sup>st</sup> step: when  $d_{ij} = \boxed{d_{12} = 1}$

$$d(3, h_{1,2}) = \frac{1}{1+1} d(1,3) + d(2,3) \frac{1}{1+1}$$

$$= \frac{11+2}{2} = 6.5$$

$$d(4, h_{1,2}) = \frac{1}{4+4} d(4,1) + d(4,2) \frac{1}{2} =$$

$$= \frac{5+3}{2} = 4$$

thus, the new dissimilarity matrix is:

$$\begin{matrix} & h_{1,2} & 3 & 4 \\ h_{1,2} & \left[ \begin{array}{ccc} 1 & 2 & 3 & 4 \\ 0 & 0 & 0 & 0 \\ 6.5 & 0 & 0 & 0 \\ 4 & 4 & 0 & 0 \end{array} \right] \\ 3 & & & \\ 4 & & & \end{matrix}$$

2<sup>nd</sup> step:

when  $d_{ij} = \boxed{d(h_{1,2}, 4)} = \boxed{d(3, 4)} = 4$

$$d(3, h_{1,2,4}) = d(3,4) \frac{1}{1+2} + d(3, h_{1,2}) \frac{2}{1+2}$$

$$= \frac{4 + 6.5 * 2}{3} = \frac{17}{3} = 5.67$$

$$D = \begin{matrix} & 3 & & \\ h_{1,2,4} & \left[ \begin{array}{ccc} 0 & & \\ 5.67 & 0 & \end{array} \right]; & \text{when } d_{ij} = 5.67 = \\ & & & = d(3, h_{1,2,4}) \end{matrix}$$

1. An observation  $x$  comes from one of the two populations with prior probabilities  $P(Y = 0) = P(Y = 1)$  and probability mass functions:

$$P(X = x|Y = j) = \begin{cases} \frac{\lambda_j^x e^{-\lambda_j}}{x!}, & x = 0, 1, 2, \dots \\ 0, & \text{otherwise,} \end{cases}$$

with  $\lambda_1 > \lambda_0 > 0$ ,  $j = 0, 1$ .

- (a) Obtain the classification rule that minimizes the total probability of misclassification. (4.0)

$$\begin{aligned} \text{TPR} &= P(\text{classifying } x \text{ with } Y=1, Y=0) + \\ &\quad + P(\text{classifying } x \text{ with } Y=0, Y=1) \\ &= \sum_{x \in R_1} P(x=x|Y=0)P(Y=0) + \\ &\quad + \sum_{x \in R_0} P(x=x|Y=1)P(Y=1) \\ &= \left[ 1 - \sum_{x \in R_0} P(x=x|Y=0) \right] P(Y=0) \\ &\quad + \sum_{x \in R_0} P(x=x|Y=1)P(Y=1) \end{aligned}$$

then TPR we define  $R_1$  as the set of  
 $x \in \mathbb{N}_0$  such that

$$P(x=x|Y=1)P(Y=1) - P(x=x|Y=0)P(Y=0) \geq 0$$

where in our case  $P(Y=1) = P(Y=0) = \frac{1}{2}$ , thus

Assign  $x_0$  to  $\gamma = 1$  iff (since  $f_1 > f_0$ )

$$P(X=x_0 | \gamma=1) \geq P(X=x_0 | \gamma=0)$$

$$\Leftrightarrow e^{-d_1} \frac{\frac{x_0}{d_1}}{x_0!} \geq e^{-d_0} \frac{\frac{x_0}{d_0}}{x_0!}$$

$$\Leftrightarrow e^{-(d_1-d_0)} \left(\frac{d_1}{d_0}\right)^{x_0} \geq 1 \quad \text{where } d_1 > d_0$$
$$\Leftrightarrow \frac{d_1}{d_0} \geq 1$$

$$\Leftrightarrow \left(\frac{d_1}{d_0}\right)^{x_0} \geq e^{(d_1-d_0)}$$

$$\Leftrightarrow x_0 \log\left(\frac{d_1}{d_0}\right) \geq d_1 - d_0$$

$$\Leftrightarrow x_0 \geq \frac{d_1 - d_0}{\log(d_1) - \log(d_0)}, \text{ since } d_1 > d_0$$

Thus, the optimal classification rule is:

$$\begin{cases} \text{If } x_0 \geq \frac{d_1 - d_0}{\log(d_1) - \log(d_0)} \Rightarrow \text{assign } x_0 \text{ to } \gamma = 1 \\ \text{otherwise} \Rightarrow \text{assign } x_0 \text{ to } \gamma = 0 \end{cases}$$

1. An observation  $x$  comes from one of the two populations with prior probabilities  $P(Y = 0) = P(Y = 1)$  and probability mass functions:

$$P(X = x|Y = j) = \begin{cases} \frac{\lambda_j^x e^{-\lambda_j}}{x!}, & x = 0, 1, 2, \dots \\ 0, & \text{otherwise,} \end{cases}$$

with  $\lambda_1 > \lambda_0 > 0$ ,  $j = 0, 1$ .

- (b) Assuming  $\lambda_0 = 1$ , and  $\lambda_1 = 10$ , obtain:

i. The recall of class  $\frac{1}{4}$  (3.0)

ii. The precision of class 1. (2.0)

iii. The  $F_1$  measure of class 1. (1.0)

If  $d_0 = 1$  and  $d_1 = 10$  then

$$m = \frac{d_1 - d_0}{\log(d_1) - \log(d_0)} = \frac{4 - 1}{\log(4) - \log(1)} \approx 2.164$$

$$\begin{aligned} \text{(i) } \text{Rec}(1) &= P(\text{classify } x \text{ in } Y=1 \mid Y=1) = \\ &= P(X \geq 2.164 \mid Y=1) = 1 - P(X \leq 2.164) \\ &= 1 - F_{P_0(4)}(2) = 1 - 0.2381 \\ &= 0.7619 \end{aligned}$$

$$\begin{aligned} \text{(ii) } \text{Pre}(1) &= P(Y=1 \mid \text{classify in } Y=1) = \\ &= \frac{P(\text{classify in } Y=1 \mid Y=1) p_1}{P(\text{classify in } Y=1 \mid Y=1) p_1 + P(\text{classify in } Y=1 \mid Y=0) p_0} \end{aligned}$$

$$\text{Pre}(1) = \frac{0.7619}{0.7619 + (1 - 0.90465)} = 0.90465$$

$$\begin{aligned}
 & P(\text{Classify } x \text{ as } y=1 \mid y=0) = \\
 & = 1 - P(x < 2.165 \mid y=0) = 1 - P(x \leq 2 \mid y=0) \\
 & = 1 - F_{P(x \mid y=0)}(2) = 1 - 0.9197 = 0.0803
 \end{aligned}$$

$$\text{(iii)} \quad F_1(1) = \frac{2 \cdot \text{Rec}(1) \cdot \text{Pr}(1)}{\text{Rec}(1) + \text{Pr}(1)} = 0.8272$$

(c) Comment your results.

(1.0)

The Precision, and  $F_1$ -measure of class  $y=1$  are very high but  $\text{Rec}(1)$  is not so high (0.7619)

thus, to be sure about the classifier's performance it is important to consider an overall measure of performance, like:

$$\begin{aligned}
 \text{Overall accuracy} &= \text{Rec}(1) p_1 + \text{Rec}(0) p_0 \\
 &= \frac{0.7619 + (1 - 0.0803)}{2} \\
 &= 0.8408
 \end{aligned}$$

is fairly high, meaning that there is a good separation between the groups leading to a high overall accuracy. nevertheless the probability of correctly assign an observation of  $y=1$  is not as high as we could like.