

Aprendizagem 2023/24  
Homework III

**I. Pen-and-paper [12v]**

Given the following observations,  $\left\{ \begin{pmatrix} 1 \\ 0.6 \\ 0.1 \end{pmatrix}, \begin{pmatrix} 0 \\ -0.4 \\ 0.8 \end{pmatrix}, \begin{pmatrix} 0 \\ 0.2 \\ 0.5 \end{pmatrix}, \begin{pmatrix} 1 \\ 0.4 \\ -0.1 \end{pmatrix} \right\}$ .

Consider a Bayesian clustering that assumes  $\{y_1\} \perp \{y_2, y_3\}$ , two clusters following a Bernoulli distribution on  $y_1$  ( $p_1$  and  $p_2$ ), a multivariate Gaussian on  $\{y_2, y_3\}$  ( $N_1$  and  $N_2$ ), and the following initial mixture:

$$\pi_1 = 0.5, \pi_2 = 0.5$$

$$p_1 = P(y_1 = 1) = 0.3, \quad p_2 = P(y_1 = 1) = 0.7$$

$$N_1 \left( \mathbf{u}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \Sigma_1 = \begin{pmatrix} 2 & 0.5 \\ 0.5 & 2 \end{pmatrix} \right), \quad N_2 \left( \mathbf{u}_2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma_2 = \begin{pmatrix} 1.5 & 1 \\ 1 & 1.5 \end{pmatrix} \right).$$

Note: you can solve this exercise by neglecting  $y_1$  and still scoring up to 70% of its grade.

1) [6v] Perform one epoch of the EM clustering algorithm and determine the new parameters. Hint: we suggest you to use numpy and scipy, however disclose the intermediary results step by step.

Homework 4:  
Clustering and PCA  
I. Pen-and-Paper:

Exercício 1)  
E-Step (Notes):

Step 1:  $P(X_n | C_k) = \begin{cases} p_1 \cdot \mathcal{N}(X_n | \mu_{k1}, \Sigma_{k1}), & y_1=1 \text{ and } k=1 \\ (1-p_1) \cdot \mathcal{N}(X_n | \mu_{k1}, \Sigma_{k1}), & y_1=0 \text{ and } k=1 \\ p_2 \cdot \mathcal{N}(X_n | \mu_{k2}, \Sigma_{k2}), & y_1=1 \text{ and } k=2 \\ (1-p_2) \cdot \mathcal{N}(X_n | \mu_{k2}, \Sigma_{k2}), & y_1=0 \text{ and } k=2 \end{cases}$

$p_1 = 0.3, p_2 = 0.7, \mathcal{N}(X_n | \mu_{k1}, \Sigma_{k1}) = \frac{1}{(2\pi)^{\frac{m}{2}} \sqrt{|\Sigma|}} e^{-\frac{1}{2}(X_n - \mu_{k1})^T \Sigma_{k1}^{-1} (X_n - \mu_{k1})}$

Step 2:  $P(C_k, X_n) = \pi_k \cdot P(X_n | C_k)$

Step 3:  $P(X_n) = \sum_{k=1}^K P(C_k, X_n), \quad \gamma(C_k | X_n) = \frac{P(C_k, X_n)}{P(X_n)}$

Auxiliary Calculations:

$|\Sigma_1| = 4 - 0.5^2 = 3.75, \quad \Sigma_1^{-1} = \frac{1}{3.75} \begin{bmatrix} 2 & -0.5 \\ -0.5 & 2 \end{bmatrix} = \begin{bmatrix} 0.53 & -0.13 \\ -0.13 & 0.53 \end{bmatrix}$

$|\Sigma_2| = 1.5^2 - 1 = 1.25, \quad \Sigma_2^{-1} = \frac{1}{1.25} \begin{bmatrix} 1.5 & -1 \\ -1 & 1.5 \end{bmatrix} = \begin{bmatrix} 1.2 & -0.8 \\ -0.8 & 1.2 \end{bmatrix}$

E-Step:

Step 1:

$P(X_1 | C_1) = P(y_1=1 | C_1) \cdot \mathcal{N} \left( \begin{bmatrix} 0.6 \\ 0.1 \end{bmatrix} \middle| \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 & 0.5 \\ 0.5 & 2 \end{bmatrix} \right)$

Using the auxiliary calculations:

$P(X_1 | C_1) = 0.3 \cdot \frac{1}{(2\pi)^{\frac{2}{2}} \sqrt{3.75}} e^{-\frac{1}{2} \left( \begin{bmatrix} 0.6 \\ 0.1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^T \begin{bmatrix} 0.53 & -0.13 \\ -0.13 & 0.53 \end{bmatrix} \left( \begin{bmatrix} 0.6 \\ 0.1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)}$

$= \frac{0.3}{2\pi \sqrt{3.75}} e^{-\frac{1}{2} \cdot [-0.4 \ -0.9] \cdot \begin{bmatrix} -0.095 \\ -0.425 \end{bmatrix}} = \frac{0.3}{2\pi \sqrt{3.75}} e^{-0.21025} \approx 0.0200$

$\Rightarrow P(X_1 | C_1) \approx 0.0200$

$P(X_2 | C_1) = P(y_1=0 | C_1) \cdot \mathcal{N} \left( \begin{bmatrix} -0.4 \\ 0.8 \end{bmatrix} \middle| \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 & 0.5 \\ 0.5 & 2 \end{bmatrix} \right)$

Using scipy.stats.multivariate\_normal:

$P(X_2 | C_1) = 0.7 \cdot 0.0500 \approx 0.0350$

$P(X_3 | C_1) = P(y_1=0 | C_1) \cdot \mathcal{N} \left( \begin{bmatrix} 0.2 \\ 0.5 \end{bmatrix} \middle| \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 & 0.5 \\ 0.5 & 2 \end{bmatrix} \right)$

Using scipy.stats.multivariate\_normal:

$P(X_3 | C_1) = 0.7 \cdot 0.0684 \approx 0.0479$

$P(X_4 | C_1) = P(y_1=1 | C_1) \cdot \mathcal{N} \left( \begin{bmatrix} 0.4 \\ -0.1 \end{bmatrix} \middle| \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 & 0.5 \\ 0.5 & 2 \end{bmatrix} \right)$

Using scipy.stats.multivariate\_normal:

$P(X_4 | C_1) = 0.3 \cdot 0.0540 \approx 0.0177$

$\rightarrow p(x_1|c_2) = p(y_1=1|c_2) \cdot \mathcal{N}\left(\begin{bmatrix} 0,6 \\ 0,1 \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1,5 & 1 \\ 1 & 1,5 \end{bmatrix}\right),$   
 Using the auxiliary calculations:  

$$p(x_1|c_2) = 0,7 \cdot \frac{1}{(2\pi)^{\frac{2}{2}} \sqrt{|1,25|}} e^{\left(-\frac{1}{2} \left(\begin{bmatrix} 0,6 \\ 0,1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix}\right)^T \begin{bmatrix} 1,2 & -0,8 \\ -0,8 & 1,2 \end{bmatrix} \left(\begin{bmatrix} 0,6 \\ 0,1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix}\right)\right)}$$
  

$$= \frac{0,7}{2\pi \sqrt{1,25}} e^{\left(-\frac{1}{2} \cdot [0,6 \ -0,1] \cdot \begin{bmatrix} 0,64 \\ -0,36 \end{bmatrix}\right)} = \frac{0,7}{2\pi \sqrt{1,25}} e^{-0,174} \Leftrightarrow$$
  
 $\Leftrightarrow p(x_1|c_2) \approx 0,0837$   
 $\cdot p(x_2|c_2) = p(y_1=0|c_2) \cdot \mathcal{N}\left(\begin{bmatrix} -0,4 \\ 0,8 \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1,5 & 1 \\ 1 & 1,5 \end{bmatrix}\right),$   
 Using scipy.stats.multivariate\_normal:  
 $p(x_2|c_2) = 0,3 \cdot 0,0682 \approx 0,0205$   
 $\cdot p(x_3|c_2) = p(y_1=0|c_2) \cdot \mathcal{N}\left(\begin{bmatrix} 0,2 \\ 0,5 \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1,5 & 1 \\ 1 & 1,5 \end{bmatrix}\right),$   
 Using scipy.stats.multivariate\_normal:  
 $p(x_3|c_2) = 0,3 \cdot 0,130 \approx 0,0389$   
 $\cdot p(x_4|c_2) = p(y_1=1|c_2) \cdot \mathcal{N}\left(\begin{bmatrix} -0,4 \\ 0,1 \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1,5 & 1 \\ 1 & 1,5 \end{bmatrix}\right),$   
 Using scipy.stats.multivariate\_normal:  
 $p(x_4|c_2) = 0,7 \cdot 0,125 \approx 0,0872$   
Step 2:  $\pi_1 = 0,5$   $\pi_2 = 0,5$   
 $\cdot p(x_1, c_1) = \pi_1 \cdot p(x_1|c_1) \approx 0,00991$   $\cdot p(x_2, c_1) = \pi_1 \cdot p(x_2|c_1) \approx 0,0175$   
 $\cdot p(x_3, c_1) = \pi_1 \cdot p(x_3|c_1) \approx 0,0239$   $\cdot p(x_4, c_1) = \pi_1 \cdot p(x_4|c_1) \approx 0,00886$

$\rightarrow$   
 $\cdot p(x_1, c_2) = \pi_2 \cdot p(x_1|c_2) \approx 0,0419$   $\cdot p(x_2, c_2) = \pi_2 \cdot p(x_2|c_2) \approx 0,0102$   
 $\cdot p(x_3, c_2) = \pi_2 \cdot p(x_3|c_2) \approx 0,0194$   $\cdot p(x_4, c_2) = \pi_2 \cdot p(x_4|c_2) \approx 0,0436$   
Step 3:  
 $\cdot p(x_1) = \sum_{k=1}^K p(c_k, x_1) = p(c_1, x_1) + p(c_2, x_1) \approx 0,0519$   
 $\cdot p(x_2) = \sum_{k=1}^K p(c_k, x_2) = p(c_1, x_2) + p(c_2, x_2) \approx 0,0277$   
 $\cdot p(x_3) = \sum_{k=1}^K p(c_k, x_3) = p(c_1, x_3) + p(c_2, x_3) \approx 0,0434$   
 $\cdot p(x_4) = \sum_{k=1}^K p(c_k, x_4) = p(c_1, x_4) + p(c_2, x_4) \approx 0,0524$   
 $\cdot \gamma(c_{11}) = p(c_1|x_1) = \frac{p(c_1, x_1)}{p(x_1)} \approx 0,193$   
 $\cdot \gamma(c_{12}) = p(c_2|x_1) = \frac{p(c_2, x_1)}{p(x_1)} \approx 0,807$   
 $\cdot \gamma(c_{21}) = p(c_1|x_2) = \frac{p(c_1, x_2)}{p(x_2)} \approx 0,631$   
 $\cdot \gamma(c_{22}) = p(c_2|x_2) = \frac{p(c_2, x_2)}{p(x_2)} \approx 0,369$   
 $\cdot \gamma(c_{31}) = p(c_1|x_3) = \frac{p(c_1, x_3)}{p(x_3)} \approx 0,552$   
 $\cdot \gamma(c_{32}) = p(c_2|x_3) = \frac{p(c_2, x_3)}{p(x_3)} \approx 0,448$

$$\gamma(C_{41}) = P(C_1 | X_4) = \frac{P(C_1, X_4)}{P(X_4)} \approx 0,169$$

$$\gamma(C_{42}) = P(C_2 | X_4) = \frac{P(C_2, X_4)}{P(X_4)} \approx 0,831$$

M-Step (Notes):

Step 1:

$$N_K = \sum_{n=1}^N \gamma(C_{nK})$$

Step 2:

$$P_K = \frac{1}{N_K} \sum_{n=1}^N \gamma(C_{nK}) \cdot X_n$$

Step 3:

$$\mu_K = \frac{1}{N_K} \sum_{n=1}^N \gamma(C_{nK}) \cdot X_n$$

Step 4:

$$\sum_K = \frac{1}{N_K} \sum_{n=1}^N \gamma(C_{nK}) \cdot (X_n - \mu_K) \cdot (X_n - \mu_K)^T$$

Step 5:

$$\pi_K = \frac{N_K}{N}$$

M-Step:

Step 1:

$$N_1 = \sum_{n=1}^4 \gamma(C_{n1}) = \gamma(C_{11}) + \gamma(C_{21}) + \gamma(C_{31}) + \gamma(C_{41}) =$$

$$= 0,193 + 0,631 + 0,552 + 0,169 \Leftrightarrow N_1 = 1,545$$

$$N_2 = \sum_{n=1}^4 \gamma(C_{n2}) = \gamma(C_{12}) + \gamma(C_{22}) + \gamma(C_{32}) + \gamma(C_{42}) =$$

$$= 0,807 + 0,369 + 0,448 + 0,831 \Leftrightarrow N_2 = 2,455$$

Step 2:

$$P_1 = \frac{1}{N_1} \sum_{n=1}^4 \gamma(C_{n1}) \cdot X_n = \frac{1}{N_1} (\gamma(C_{11}) \cdot X_1 + \gamma(C_{21}) \cdot X_2 + \gamma(C_{31}) \cdot X_3 + \gamma(C_{41}) \cdot X_4) =$$

$$= \frac{0,193 \cdot 1 + 0,631 \cdot 0 + 0,552 \cdot 0 + 0,169 \cdot 1}{1,545} \approx 0,234$$

$$P_2 = \frac{1}{N_2} \sum_{n=1}^4 \gamma(C_{n2}) \cdot X_n = \frac{1}{N_2} (\gamma(C_{12}) \cdot X_1 + \gamma(C_{22}) \cdot X_2 + \gamma(C_{32}) \cdot X_3 + \gamma(C_{42}) \cdot X_4) =$$

$$= \frac{0,807 \cdot 1 + 0,369 \cdot 0 + 0,448 \cdot 0 + 0,831 \cdot 1}{2,455} \approx 0,667$$

Step 3:

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^4 \gamma(C_{n1}) \cdot X_n = \frac{1}{N_1} (\gamma(C_{11}) \cdot X_1 + \gamma(C_{21}) \cdot X_2 + \gamma(C_{31}) \cdot X_3 + \gamma(C_{41}) \cdot X_4) =$$

$$= \frac{0,193 \begin{bmatrix} 0,67 \\ 0,1 \end{bmatrix} + 0,631 \begin{bmatrix} -0,4 \\ 0,8 \end{bmatrix} + 0,552 \begin{bmatrix} 0,27 \\ 0,5 \end{bmatrix} + 0,169 \begin{bmatrix} 0,47 \\ -0,1 \end{bmatrix}}{1,545} =$$

$$= \frac{1}{1,545} \begin{bmatrix} 0,0414 \\ 0,7832 \end{bmatrix} = \begin{bmatrix} 0,0268 \\ 0,5069 \end{bmatrix}$$

$$\mu_2 = \frac{1}{N_2} \sum_{n=1}^4 \gamma(C_{n2}) \cdot X_n = \frac{1}{N_2} (\gamma(C_{12}) \cdot X_1 + \gamma(C_{22}) \cdot X_2 + \gamma(C_{32}) \cdot X_3 + \gamma(C_{42}) \cdot X_4) =$$

$$= \frac{0,807 \begin{bmatrix} 0,67 \\ 0,1 \end{bmatrix} + 0,369 \begin{bmatrix} -0,4 \\ 0,8 \end{bmatrix} + 0,448 \begin{bmatrix} 0,27 \\ 0,5 \end{bmatrix} + 0,831 \begin{bmatrix} 0,47 \\ -0,1 \end{bmatrix}}{2,455} =$$

$$= \frac{1}{2,455} \begin{bmatrix} 0,7586 \\ 0,5168 \end{bmatrix} = \begin{bmatrix} 0,3091 \\ 0,2105 \end{bmatrix}$$

→ Step 4:

$$\sum_1 = \frac{1}{N_1} \sum_{n=1}^4 \gamma(C_{n1}) (X_n - \mu_1) (X_n - \mu_1)^T =$$

$$= \frac{1}{1,545} \left( 0,193 \begin{bmatrix} 0,5732 \\ -0,4096 \end{bmatrix} \begin{bmatrix} 0,5732 & -0,4096 \end{bmatrix} + 0,631 \begin{bmatrix} -0,4268 \\ 0,2904 \end{bmatrix} \begin{bmatrix} -0,4268 & 0,2904 \end{bmatrix} + \right.$$

$$\left. + 0,552 \begin{bmatrix} 0,1732 \\ -0,0096 \end{bmatrix} \begin{bmatrix} 0,1732 & -0,0096 \end{bmatrix} + 0,169 \begin{bmatrix} 0,3732 \\ -0,6096 \end{bmatrix} \begin{bmatrix} 0,3732 & -0,6096 \end{bmatrix} \right) =$$

$$= \frac{1}{1,545} \left( \begin{bmatrix} 0,0634 & -0,0453 \\ -0,0453 & 0,0324 \end{bmatrix} + \begin{bmatrix} 0,115 & -0,0782 \\ -0,0782 & 0,0532 \end{bmatrix} + \begin{bmatrix} 0,0166 & -0,00118 \\ -0,00118 & 0,00051 \end{bmatrix} + \right.$$

$$\left. + \begin{bmatrix} 0,0235 & -0,0384 \\ -0,0384 & 0,0628 \end{bmatrix} \right) = \begin{bmatrix} 0,141 & -0,105 \\ -0,105 & 0,0961 \end{bmatrix}$$

· Step 5:

$$\sum_2 = \frac{1}{N_2} \sum_{n=1}^4 \gamma(C_{n2}) (X_n - \mu_2) (X_n - \mu_2)^T =$$

$$= \frac{1}{2,455} \left( 0,807 \begin{bmatrix} 0,291 \\ -0,1105 \end{bmatrix} \begin{bmatrix} 0,291 & -0,1105 \end{bmatrix} + 0,369 \begin{bmatrix} -0,209 \\ 0,5895 \end{bmatrix} \begin{bmatrix} -0,209 & 0,5895 \end{bmatrix} + \right.$$

$$\left. + 0,448 \begin{bmatrix} -0,109 \\ 0,2895 \end{bmatrix} \begin{bmatrix} -0,109 & 0,2895 \end{bmatrix} + 0,831 \begin{bmatrix} 0,091 \\ -0,3105 \end{bmatrix} \begin{bmatrix} 0,091 & -0,3105 \end{bmatrix} \right) =$$

$$= \frac{1}{2,455} \left( \begin{bmatrix} 0,0683 & -0,0259 \\ -0,0259 & 0,00985 \end{bmatrix} + \begin{bmatrix} 0,185 & -0,154 \\ -0,154 & 0,128 \end{bmatrix} + \begin{bmatrix} 0,00532 & -0,0141 \\ -0,0141 & 0,0375 \end{bmatrix} + \right.$$

$$\left. + \begin{bmatrix} 0,00688 & -0,0235 \\ -0,0235 & 0,0801 \end{bmatrix} \right) = \begin{bmatrix} 0,108 & -0,0887 \\ -0,0887 & 0,104 \end{bmatrix}$$

$$\pi_1 = \frac{N_1}{N} = \frac{1,545}{1,545 + 2,455} \approx 0,386$$

$$\pi_2 = \frac{N_2}{N} = \frac{2,455}{1,545 + 2,455} \approx 0,614$$



- 2) [2v] Given the new observation,  $\mathbf{x}_{\text{new}} = \begin{pmatrix} 1 \\ 0.3 \\ 0.7 \end{pmatrix}$ , determine the cluster memberships (posteriors).

Exercício 2)

$$p(C=k | \mathbf{x}_{\text{new}}) = \frac{p(C=k, \mathbf{x}_{\text{new}})}{p(\mathbf{x}_{\text{new}})} = \frac{\pi_k \cdot p(\mathbf{x}_{\text{new}} | C_k)}{p(\mathbf{x}_{\text{new}})}$$

$p(\mathbf{x}_{\text{new}} | C_k) = p_k \cdot \mathcal{N}(\mathbf{x}_{\text{new}} | \mu_k, \Sigma_k) \cdot \gamma_k = 1$

For  $k=1$ :

$$p(\mathbf{x}_{\text{new}} | C_1) = p_1 \cdot \mathcal{N}(\mathbf{x}_{\text{new}} | \mu_1, \Sigma_1)$$

From the previous exercise we know that:

$$p_1 = \frac{0,234}{0,386} \quad \mu_1 = \begin{bmatrix} 0,0268 \\ 0,5069 \end{bmatrix} \quad \Sigma_1 = \begin{bmatrix} 0,141 & -0,105 \\ -0,105 & 0,0961 \end{bmatrix}$$

Therefore:

$$p(\mathbf{x}_{\text{new}} | C_1) = \frac{0,234}{0,386} \cdot \mathcal{N}\left(\begin{bmatrix} 0,3 \\ 0,7 \end{bmatrix} \middle| \begin{bmatrix} 0,0268 \\ 0,5069 \end{bmatrix}, \begin{bmatrix} 0,141 & -0,105 \\ -0,105 & 0,0961 \end{bmatrix}\right)$$

Using scipy.stats.multivariate\_normal:

$$p(\mathbf{x}_{\text{new}} | C_1) = \frac{0,234}{0,386} \cdot 0,0295 \approx 0,00690$$

For  $k=2$ :

$$p(\mathbf{x}_{\text{new}} | C_2) = p_2 \cdot \mathcal{N}(\mathbf{x}_{\text{new}} | \mu_2, \Sigma_2)$$

From the previous exercise we know that:

$$p_2 = 0,667 \quad \mu_2 = \begin{bmatrix} 0,309 \\ 0,2105 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 0,108 & -0,0887 \\ -0,0887 & 0,104 \end{bmatrix}$$

→ Therefore:

$$p(\mathbf{x}_{\text{new}} | C_2) = 0,667 \cdot \mathcal{N}\left(\begin{bmatrix} 0,3 \\ 0,7 \end{bmatrix} \middle| \begin{bmatrix} 0,309 \\ 0,2105 \end{bmatrix}, \begin{bmatrix} 0,108 & -0,0887 \\ -0,0887 & 0,104 \end{bmatrix}\right)$$

Using scipy.stats.multivariate\_normal:

$$p(\mathbf{x}_{\text{new}} | C_2) = 0,667 \cdot 0,0629 \approx 0,0417$$

$$p(\mathbf{x}_{\text{new}}) = \sum_{k=1}^K p(C_k, \mathbf{x}_{\text{new}}) = \sum_{k=1}^K \pi_k \cdot p(\mathbf{x}_{\text{new}} | C_k) =$$

$$= \pi_1 \cdot p(\mathbf{x}_{\text{new}} | C_1) + \pi_2 \cdot p(\mathbf{x}_{\text{new}} | C_2)$$

From the previous exercise we know that:

$$\pi_1 = 0,386 \quad \pi_2 = 0,614 \quad \text{Therefore:}$$

$$p(\mathbf{x}_{\text{new}}) = 0,386 \cdot 0,00690 + 0,614 \cdot 0,0417 \approx 0,0282$$

Therefore:

$$p(C_1 | \mathbf{x}_{\text{new}}) = \frac{\pi_1 \cdot p(\mathbf{x}_{\text{new}} | C_1)}{p(\mathbf{x}_{\text{new}})} = \frac{0,386 \cdot 0,00690}{0,0282} \approx 0,0943$$

$$p(C_2 | \mathbf{x}_{\text{new}}) = \frac{\pi_2 \cdot p(\mathbf{x}_{\text{new}} | C_2)}{p(\mathbf{x}_{\text{new}})} = \frac{0,614 \cdot 0,0417}{0,0282} \approx 0,906$$

3) [2.5v] Performing a hard assignment of observations to clusters under a ML assumption, identify the silhouette of the larger cluster under a Manhattan distance.

Exercício 3)  
 For the following calculations we will be using  
 scipy.stats.multivariate\_normal.

$$\begin{aligned} \mathcal{N}(X_1|C_1) &= \mathcal{N}\left(\begin{bmatrix} 0,6 \\ 0,1 \end{bmatrix} \middle| \begin{bmatrix} 0,0268 & 0,141 \\ 0,5069 & -0,105 \end{bmatrix}, \begin{bmatrix} 0,141 & -0,105 \\ 0,105 & 0,0961 \end{bmatrix}\right) \approx 0,980 \\ \mathcal{N}(X_2|C_1) &= \mathcal{N}\left(\begin{bmatrix} -0,4 \\ 0,8 \end{bmatrix} \middle| \begin{bmatrix} 0,0268 & 0,141 \\ 0,5069 & -0,105 \end{bmatrix}, \begin{bmatrix} 0,141 & -0,105 \\ 0,105 & 0,0961 \end{bmatrix}\right) \approx 1,636 \\ \mathcal{N}(X_3|C_1) &= \mathcal{N}\left(\begin{bmatrix} 0,2 \\ 0,5 \end{bmatrix} \middle| \begin{bmatrix} 0,0268 & 0,141 \\ 0,5069 & -0,105 \end{bmatrix}, \begin{bmatrix} 0,141 & -0,105 \\ 0,105 & 0,0961 \end{bmatrix}\right) \approx 1,879 \\ \mathcal{N}(X_4|C_1) &= \mathcal{N}\left(\begin{bmatrix} 0,4 \\ -0,1 \end{bmatrix} \middle| \begin{bmatrix} 0,0268 & 0,141 \\ 0,5069 & -0,105 \end{bmatrix}, \begin{bmatrix} 0,141 & -0,105 \\ 0,105 & 0,0961 \end{bmatrix}\right) \approx 0,0928 \\ \mathcal{N}(X_1|C_2) &= \mathcal{N}\left(\begin{bmatrix} 0,6 \\ 0,1 \end{bmatrix} \middle| \begin{bmatrix} 0,309 & 0,108 \\ 0,2105 & -0,0887 \end{bmatrix}, \begin{bmatrix} 0,108 & -0,0887 \\ 0,104 & -0,0887 \end{bmatrix}\right) \approx 1,423 \\ \mathcal{N}(X_2|C_2) &= \mathcal{N}\left(\begin{bmatrix} -0,4 \\ 0,8 \end{bmatrix} \middle| \begin{bmatrix} 0,309 & 0,108 \\ 0,2105 & -0,0887 \end{bmatrix}, \begin{bmatrix} 0,108 & -0,0887 \\ 0,104 & -0,0887 \end{bmatrix}\right) \approx 0,270 \\ \mathcal{N}(X_3|C_2) &= \mathcal{N}\left(\begin{bmatrix} 0,2 \\ 0,5 \end{bmatrix} \middle| \begin{bmatrix} 0,309 & 0,108 \\ 0,2105 & -0,0887 \end{bmatrix}, \begin{bmatrix} 0,108 & -0,0887 \\ 0,104 & -0,0887 \end{bmatrix}\right) \approx 1,366 \\ \mathcal{N}(X_4|C_2) &= \mathcal{N}\left(\begin{bmatrix} 0,4 \\ -0,1 \end{bmatrix} \middle| \begin{bmatrix} 0,309 & 0,108 \\ 0,2105 & -0,0887 \end{bmatrix}, \begin{bmatrix} 0,108 & -0,0887 \\ 0,104 & -0,0887 \end{bmatrix}\right) \approx 1,0773 \end{aligned}$$

Posteriors:

$$\begin{aligned} P(X_1|C_1) &= P(Y_1=1|C_1) \cdot \mathcal{N}(X_1|C_1) = 0,234 \times 0,980 \approx 0,229 \\ P(X_2|C_1) &= P(Y_1=0|C_1) \cdot \mathcal{N}(X_2|C_1) = 0,766 \times 1,636 \approx 1,254 \\ P(X_3|C_1) &= P(Y_1=0|C_1) \cdot \mathcal{N}(X_3|C_1) = 0,766 \times 1,879 \approx 1,439 \\ P(X_4|C_1) &= P(Y_1=1|C_1) \cdot \mathcal{N}(X_4|C_1) = 0,234 \times 0,0928 \approx 0,0217 \\ P(X_1|C_2) &= P(Y_1=1|C_2) \cdot \mathcal{N}(X_1|C_2) = 0,667 \times 1,423 \approx 0,949 \\ P(X_2|C_2) &= P(Y_1=0|C_2) \cdot \mathcal{N}(X_2|C_2) = 0,333 \times 0,270 \approx 0,0898 \\ P(X_3|C_2) &= P(Y_1=0|C_2) \cdot \mathcal{N}(X_3|C_2) = 0,333 \times 1,366 \approx 0,455 \\ P(X_4|C_2) &= P(Y_1=1|C_2) \cdot \mathcal{N}(X_4|C_2) = 0,667 \times 1,0773 \approx 0,719 \end{aligned}$$

→ Under ML assumption we want:

$$\begin{aligned} \arg \max_{W \in \{C_1, C_2\}} P(X_1|W) &= C_2 \\ \arg \max_{W \in \{C_1, C_2\}} P(X_2|W) &= C_1 \\ \arg \max_{W \in \{C_1, C_2\}} P(X_3|W) &= C_1 \\ \arg \max_{W \in \{C_1, C_2\}} P(X_4|W) &= C_2 \end{aligned}$$

clusters =  $\{C_1 = \{X_2, X_3\}, C_2 = \{X_1, X_4\}\}$   
 So both clusters have the same size.

$$\begin{aligned} a(X_1) &= |X_1 - X_4| = 0,4 \\ b(X_1) &= \frac{|X_1 - X_2| + |X_1 - X_3|}{2} = \frac{(2,7 + 1,8)}{2} = 2,25 \\ a(X_2) &= |X_2 - X_3| = 0,9 \\ b(X_2) &= \frac{|X_2 - X_1| + |X_2 - X_4|}{2} = \frac{(2,7 + 2,7)}{2} = 2,7 \\ a(X_3) &= |X_3 - X_2| = 0,9 \\ b(X_3) &= \frac{|X_3 - X_1| + |X_3 - X_4|}{2} = \frac{(1,8 + 1,8)}{2} = 1,8 \\ a(X_4) &= |X_4 - X_1| = 0,4 \\ b(X_4) &= \frac{|X_4 - X_2| + |X_4 - X_3|}{2} = \frac{(2,7 + 1,8)}{2} = 2,25 \end{aligned}$$

$$\begin{aligned} S(X_1) &= 1 - \frac{a(X_1)}{b(X_1)} = 0,822 & S(X_3) &= 1 - \frac{a(X_3)}{b(X_3)} = 0,5 \\ S(X_2) &= 1 - \frac{a(X_2)}{b(X_2)} = 0,6667 & S(X_4) &= 1 - \frac{a(X_4)}{b(X_4)} = 0,822 \\ S(C_1) &= \frac{S(X_2) + S(X_3)}{2} = 0,58 & S(C_2) &= \frac{S(X_1) + S(X_4)}{2} = 0,822 \end{aligned}$$

4) [0.5v] Knowing the purity of the clustering solution is 0.75, identify the number of possible classes (ground truth).

Exercício 4)

Knowing that purity = 0.75 and the formula is  $\text{purity} = \frac{1}{n} \sum_{k=1}^K \max_j (|C_k \cap L_j|)$ . For this statement,  $n=4$  (number of observations) and there are 2 clusters ( $C_1$  and  $C_2$ ):

$$0.75 = \frac{1}{4} \sum_{k=1}^2 \max_j (|C_k \cap L_j|) \Leftrightarrow$$

$$\Leftrightarrow \sum_{k=1}^2 \max_j (|C_k \cap L_j|) = 3$$

In exercise 3 it is concluded that  $C_1 = \{x_2, x_3\}$  and  $C_2 = \{x_1, x_4\}$ , i.e., a possible distribution by classes so that the sum above is equal to 3 would be  $x_2$  and  $x_3$  being of the same class and  $x_1$  and  $x_4$  being of different classes, so  $\sum_{k=1}^2 \max_j (|C_k \cap L_j|) = \max(2,0) + \max(1,1) = 2+1=3$ . Therefore a possible value for the number of classes (ground truth) is 2.

## Aprendizagem 2023/24

### Homework III

#### Programming and critical analysis [8v]

Recall the `column_diagnosis.arff` dataset from previous homework's. For the following exercises, normalize the data using sklearn's `MinMaxScaler`.

1) [4v] Using sklearn, apply k-means clustering fully unsupervised Ly on the normalized data with  $k \in \{2,3,4,5\}$  (random=0 and remaining parameters as default). Assess the silhouette and purity of the produced solutions.

```
import pandas as pd
import numpy as np
from scipy.io.arff import loadarff
from sklearn.preprocessing import MinMaxScaler
from sklearn.cluster import KMeans
from sklearn.metrics.cluster import contingency_matrix
from sklearn.metrics import silhouette_score

# Read ARFF file
data = loadarff('column_diagnosis.arff')
df = pd.DataFrame(data[0])
df['class'] = df['class'].str.decode('utf-8')

X = df.drop('class', axis=1)
y = df['class']
X_scaled = MinMaxScaler().fit_transform(X)

k_values = [2, 3, 4, 5]

for k in k_values:
    # Perform K-means clustering
    kmeans = KMeans(n_clusters=k, random_state=0).fit(X_scaled)
    if (k==3):
        kmeans3 = kmeans

    # Calculate silhouette scores and purity scores
    s_score = silhouette_score(X_scaled, kmeans.labels_)
    confusion_matrix = contingency_matrix(y, kmeans.labels_)
    p_score = np.sum(np.amax(confusion_matrix, axis=0)) / np.sum(confusion_matrix)

    print('For K =', k)
    print('Silhouette score:', s_score)
    print('Purity score:', p_score, '\n')
```

Output:

For K = 2

Silhouette score: 0.36044124340441114

Purity score: 0.632258064516129

For K = 3

Silhouette score: 0.29579055730002257

Purity score: 0.667741935483871

For K = 4

Silhouette score: 0.27442402122340176

Purity score: 0.6612903225806451

For K = 5

Silhouette score: 0.23823928397844843

Purity score: 0.6774193548387096

2) [2v] Consider the application of PCA after the data normalization:



## Aprendizagem 2023/24

### Homework III

i. Identify the variability explained by the top two principal components.

```
from sklearn.decomposition import PCA

pca = PCA(n_components=2)
principal_components = pca.fit(X_scaled)

# Compute the explained variance for the top two components
explained_variance = pca.explained_variance_ratio_
print('Explained variance of the top two components:', explained_variance)
```

Output:

Explained variance of the top two components: [0.56181445 0.20955953]

ii. For each one of these two components, sort the input variables by relevance by inspecting the absolute weights of the linear projection.

```
# Get the absolute weights
component_weights = abs(pca.components_)

# Sort the input variables by relevance
relevance_df = pd.DataFrame(component_weights, columns=X.columns, index=['Component 1', 'Component 2']).T
relevance_df = relevance_df.sort_values(by=['Component 1', 'Component 2'], ascending=False)

print('Sorted relevance of input variables in the top two components:')
print(relevance_df)
```

Output:

Sorted relevance of input variables in the top two components:

	Component 1	Component 2
pelvic_incidence	0.591621	0.100037
lumbar_lordosis_angle	0.515085	0.080047
pelvic_tilt	0.467039	0.670373
sacral_slope	0.325689	0.443303
degree_spondylolisthesis	0.216930	0.004583
pelvic_radius	0.115824	0.581074

## Aprendizagem 2023/24

### Homework III

3) [2v] Visualize side-by-side the data using: i) the ground diagnoses, and ii) the *previously* learned  $k = 3$  clustering solution. To this end, projected the normalized data onto a 2-dimensional data space using PCA and then color observations using the reference and cluster annotations.

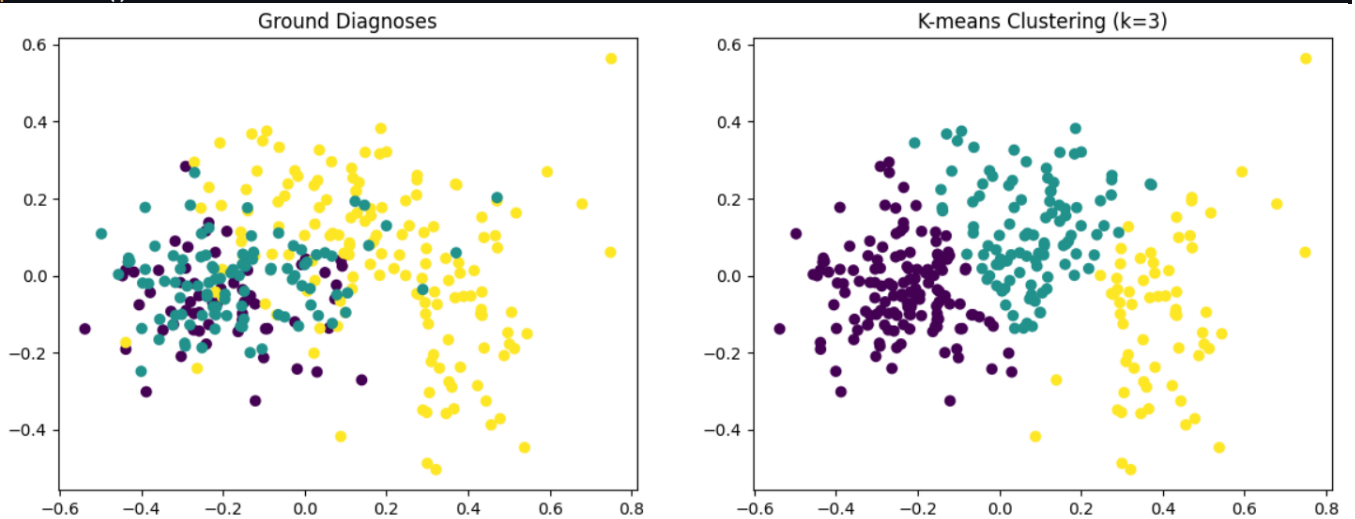
```
import matplotlib.pyplot as plt

# Reduce the dimensionality of data
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X_scaled)

# Real labels as integers to compare with kmeans labels
codes = {'Hernia':0, 'Normal':1, 'Spondylolisthesis':2}
y_pred = y.map(codes).tolist()

plt.figure(figsize=(14, 5))
plt.subplot(121)
plt.scatter(X_pca[:,0], X_pca[:,1], c=y_pred)
plt.title('Ground Diagnoses')

plt.subplot(122)
plt.scatter(X_pca[:,0], X_pca[:,1], c=kmeans3.labels_)
plt.title('K-means Clustering (k=3)')
plt.show()
```



4) [1v] Considering the results from questions (1) and (3), identify two ways on how clustering can be used to characterize the population of ill and healthy individuals.

Answer: Clustering can be used to identify groups of patients with similar medical characteristics to discover potential causes of the disease or define groups at higher risk of developing the disease.

**END**