

I. Pen-and-paper

1) i. $\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$; $\Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$; $\pi_1 = P(C_1=1) = 0,7$; $\pi_2 = P(C_2=1) = 0,3$
 $K=2$
 $\mu_1 = x_1 = \begin{pmatrix} 2 \\ 4 \end{pmatrix}$ $\mu_2 = x_2 = \begin{pmatrix} -1 \\ -4 \end{pmatrix}$

• E - skip

$C_1=1$:

$$P(x_i | C_1=1) = N(x_i | \mu_1, \Sigma_1) = \frac{1}{(2\pi)^{\frac{K}{2}}} \cdot |\Sigma_1|^{-\frac{1}{2}} \cdot \exp\left(-\frac{1}{2}(x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1)\right)$$

$$\begin{aligned} P(x_1 | C_1=1) &= 0,1592 & P(x_1, C_1) &= P(x_1 | C_1=1) \cdot \pi_1 = 0,1114 \\ P(x_2 | C_1=1) &= 2,2381 \times 10^{-17} & P(x_2, C_1) &= P(x_2 | C_1=1) \cdot \pi_1 = 1,5674 \times 10^{-17} \\ P(x_3 | C_1=1) &= 0,000393 & P(x_3, C_1) &= P(x_3 | C_1=1) \cdot \pi_1 = 0,000275 \\ P(x_4 | C_1=1) &= 2,2256 \times 10^{-6} & P(x_4, C_1) &= P(x_4 | C_1=1) \cdot \pi_1 = 1,5579 \times 10^{-6} \end{aligned}$$

$$\begin{aligned} C_2: & P(x_1 | C_2=1) = 9,4388 \times 10^{-10} & P(x_1, C_2) &= P(x_1 | C_2=1) \cdot \pi_2 = 2,8316 \times 10^{-10} \\ & P(x_2 | C_2=1) = 0,07958 & P(x_2, C_2) &= P(x_2 | C_2=1) \cdot \pi_2 = 0,023947 \\ & P(x_3 | C_2=1) = 9,8206 \times 10^{-6} & P(x_3, C_2) &= P(x_3 | C_2=1) \cdot \pi_2 = 2,9462 \times 10^{-6} \\ & P(x_4 | C_2=1) = 2,9137 \times 10^{-6} & P(x_4, C_2) &= P(x_4 | C_2=1) \cdot \pi_2 = 8,4410 \times 10^{-7} \end{aligned}$$

$$P(x_1) = P(x_1, C_1=1) + P(x_1, C_2=1) = 0,1114$$

$$P(x_2) = 0,023947$$

$$P(x_3) = 1,7045 \times 10^{-4}$$

$$P(x_4) = 5,902 \times 10^{-6}$$

• Maximização

$$N_1 = \sum_{n=1}^N P(C_1=1 | x_n) = 2,8333$$

$$N_2 = \sum_{n=1}^N P(C_2=1 | x_n) = 1,1603$$

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N P(C_k=1 | x_n) \cdot x_n$$

$$\mu_1 = \frac{1}{2,8333} \left(1 \begin{pmatrix} 2 \\ 4 \end{pmatrix} + 0 \begin{pmatrix} -1 \\ -4 \end{pmatrix} + 0,9527 \begin{pmatrix} -1 \\ 2 \end{pmatrix} + 0,3530 \begin{pmatrix} 4 \\ 0 \end{pmatrix} \right) = \begin{pmatrix} 1,5654 \\ 2,1007 \end{pmatrix}$$

$$\mu_2 = \frac{1}{1,1603} \left(0 \begin{pmatrix} 2 \\ 4 \end{pmatrix} + 1 \begin{pmatrix} -1 \\ -4 \end{pmatrix} + 0,0173 \begin{pmatrix} -1 \\ 2 \end{pmatrix} + 0,1420 \begin{pmatrix} 4 \\ 0 \end{pmatrix} \right) = \begin{pmatrix} -0,3838 \\ -2,4176 \end{pmatrix}$$

Posterioris:

$$P(C_1=1 | x_1) = \frac{P(C_1=1, x_1)}{P(x_1)} = 1$$

$$P(C_2=1 | x_1) = 0$$

$$P(C_1=1 | x_2) = 0$$

$$P(C_2=1 | x_2) = 1$$

$$P(C_1=1 | x_3) = 0,9822$$

$$P(C_2=1 | x_3) = 0,0173$$

$$P(C_1=1 | x_4) = 0,8530$$

$$P(C_2=1 | x_4) = 0,1470$$

$$\Sigma_1 = \begin{pmatrix} 4,1328 & -1,1634 \\ -1,1634 & 2,6056 \end{pmatrix}$$

$$\Sigma_2 = \begin{pmatrix} 2,3014 & 2,1060 \\ 2,1060 & 2,1694 \end{pmatrix}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N P(C_k=1 | x_n) \cdot (x_n - \mu_k) \cdot (x_n - \mu_k)^T$$

$$(x_1 - \mu_1) (x_1 - \mu_1)^T = \begin{pmatrix} 0,4388 & 0,5334 \\ 0,5334 & 2,6012 \end{pmatrix}$$

$$(x_2 - \mu_1) (x_2 - \mu_1)^T = \begin{pmatrix} 6,5813 & 12,6097 \\ 12,6097 & 39,2195 \end{pmatrix}$$

$$(x_3 - \mu_1) (x_3 - \mu_1)^T = \begin{pmatrix} 0,8513 & 0,0513 \\ 0,0513 & 0,0101 \end{pmatrix}$$

$$(x_4 - \mu_1) (x_4 - \mu_1)^T = \begin{pmatrix} 5,9235 & -5,1144 \\ -5,1144 & 4,4125 \end{pmatrix}$$

$$(x_1 - \mu_2) (x_1 - \mu_2)^T = \begin{pmatrix} 5,6825 & 17,6821 \\ 17,6821 & 57,0109 \end{pmatrix}$$

$$(x_2 - \mu_2) (x_2 - \mu_2)^T = \begin{pmatrix} 0,3191 & 0,3589 \\ 0,3589 & 0,3302 \end{pmatrix}$$

$$(x_3 - \mu_2) (x_3 - \mu_2)^T = \begin{pmatrix} 0,2787 & -2,3283 \\ -2,3283 & 2,9204 \end{pmatrix}$$

$$(x_4 - \mu_2) (x_4 - \mu_2)^T = \begin{pmatrix} 0,2137 & 14,9821 \\ 14,9821 & 11,4800 \end{pmatrix}$$

$$\pi_1 = P(C_1=1) = \frac{N_1}{N} = \frac{2,8333}{4} = 0,7083$$

$$\pi_2 = P(C_2=1) = \frac{N_2}{N} = \frac{1,1603}{4} = 0,2901$$

$$x_1 - \mu_1 = \begin{pmatrix} 2 \\ 4 \end{pmatrix} - \begin{pmatrix} 1,5654 \\ 2,1007 \end{pmatrix} = \begin{pmatrix} 0,4346 \\ 1,8993 \end{pmatrix}$$

$$x_2 - \mu_1 = \begin{pmatrix} -1 \\ -4 \end{pmatrix} - \begin{pmatrix} 1,5654 \\ 2,1007 \end{pmatrix} = \begin{pmatrix} -2,5654 \\ -6,1007 \end{pmatrix}$$

$$x_3 - \mu_1 = \begin{pmatrix} -1 \\ 2 \end{pmatrix} - \begin{pmatrix} 1,5654 \\ 2,1007 \end{pmatrix} = \begin{pmatrix} -2,5654 \\ -0,1007 \end{pmatrix}$$

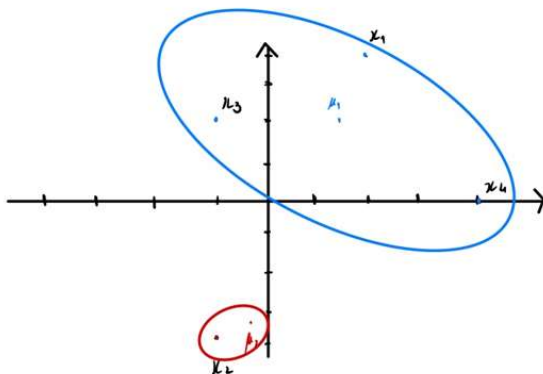
$$x_4 - \mu_1 = \begin{pmatrix} 4 \\ 0 \end{pmatrix} - \begin{pmatrix} 1,5654 \\ 2,1007 \end{pmatrix} = \begin{pmatrix} 2,4346 \\ -2,1007 \end{pmatrix}$$

$$x_1 - \mu_2 = \begin{pmatrix} 2 \\ 4 \end{pmatrix} - \begin{pmatrix} -0,3838 \\ -2,4176 \end{pmatrix} = \begin{pmatrix} 2,3838 \\ 6,4176 \end{pmatrix}$$

$$x_2 - \mu_2 = \begin{pmatrix} -1 \\ -4 \end{pmatrix} - \begin{pmatrix} -0,3838 \\ -2,4176 \end{pmatrix} = \begin{pmatrix} -0,6162 \\ -1,5824 \end{pmatrix}$$

$$x_3 - \mu_2 = \begin{pmatrix} -1 \\ 2 \end{pmatrix} - \begin{pmatrix} -0,3838 \\ -2,4176 \end{pmatrix} = \begin{pmatrix} -0,6162 \\ 4,4176 \end{pmatrix}$$

$$x_4 - \mu_2 = \begin{pmatrix} 4 \\ 0 \end{pmatrix} - \begin{pmatrix} -0,3838 \\ -2,4176 \end{pmatrix} = \begin{pmatrix} 4,3838 \\ 2,4176 \end{pmatrix}$$



Aprendizagem 2021/22
Homework IV – Group 17

2) $s = 1 - \frac{a}{b}$

$$a_1 = \frac{\|x_1 - x_3\|_2 + \|x_1 - x_4\|_2}{2} = \frac{3,6056 + 4,4721}{2} = 4,0389$$

$$s_1 = 1 - \frac{a_1}{b_1} = 0,5273$$

$$s_2 = 0, \text{ por definição pois } \|C_2\| = 1.$$

$$a_3 = \frac{\|x_3 - x_1\|_2 + \|x_3 - x_4\|_2}{2} = \frac{3,6056 + 5,3852}{2} = 4,4954$$

$$s_3 = 1 - \frac{a_3}{b_3} = 0,2507$$

$$s_4 = 1 - \frac{a_4}{b_4} = 0,2303$$

$$a_4 = \frac{\|x_4 - x_1\|_2 + \|x_4 - x_3\|_2}{2} = \frac{4,4721 + 5,3852}{2} = 4,9287$$

$$b_1 = \|x_1 - x_2\|_2 = 8,5440$$

$$s_{C_1} = 0; \quad s_{C_1} = \frac{s_1 + s_3 + s_4}{3} = \frac{0,5273 + 0,2507 + 0,2303}{3} = 0,3361$$

$$b_3 = \|x_3 - x_2\|_2 = 6$$

$$b_4 = \|x_4 - x_2\|_2 = 6,4031$$

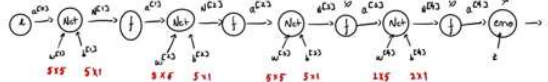
$$S = \frac{0 + 0,3361}{2} = 0,1681$$

Verifica-se que a silhueta obtida é bastante baixa (muito inferior a 1), logo seriam necessários mais clusters para melhor classificação dos pontos.

3)

a) 3)

a) MLP, inner layers: (5, 5, 5)



$$VC\text{-Dimension} = 3 \times 5^2 + 3 \times 5 + 2 \times 5 + 1 \times 1 = 101$$

$$3 \times n^2 + 3 \times n + 2 \times n + 1$$

Decision tree, 3 bins:

$$VC\text{-Dimension} = 3^3 = 3^5 = 243$$

$$3^n$$

Bayes Classifier, multivariate Gauss likelihood

$$\mu \rightarrow 5 \text{ parâmetros}$$

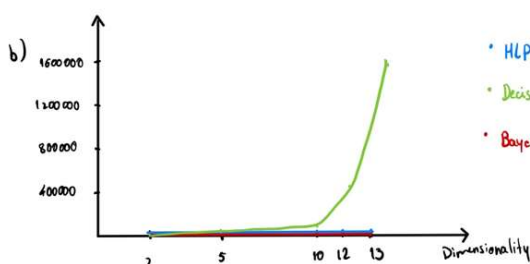
$$\Sigma \rightarrow \frac{5^2 - 5}{2} + 5 = 15$$

Por classe $P(C) \rightarrow 1 \text{ parâmetro (apenas 2 classes)}$

$$\left(\frac{n^2 + n}{2} + n \right) \times 2 + 1$$

$$VC\text{-Dimension} = (15 + 5) \times 2 + 1 = 41$$

b)



* HLP $\{ (2, 24), (5, 102), (10, 352), (12, 494), (13, 574) \}$

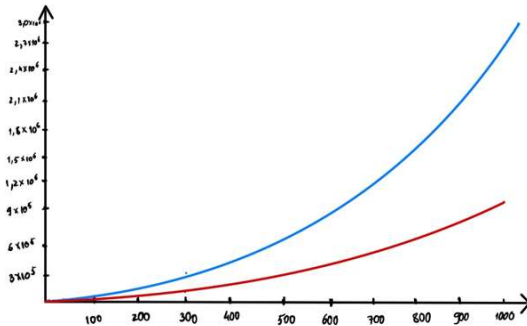
* Decision Tree $\{ (2, 9), (5, 243), (10, 59043), (12, 631441), (13, 1594323) \}$

* Bayesian classification $\{ (2, 11), (5, 41), (10, 131), (12, 161), (13, 209) \}$

É possível concluir que com o aumento da dimensionalidade da amostra, a VC-Dimension da decision tree tem um crescimento muito mais acentuado que tanto o HLP e o Bayes Classifier. ($O(3^2) \gg O(2^2)$)

c)

c)



* MLP $\{(1, 24), (5, 102), (10, 352), (30, 2352), (100, 30502), (300, 271502), (1000, 8005002)\}$

* Bayesian classification $\{(1, 11), (5, 41), (10, 151), (30, 991), (100, 10301), (300, 90901), (1000, 1003001)\}$

Apesar de ambas os modelos terem VC-Dimension de crescimento quadrático,
o MLP tem maior crescimento e consequentemente maior VC-Dimension.

II. Programming and critical analysis

4)

a) $k = 2$: ECR = 13.5

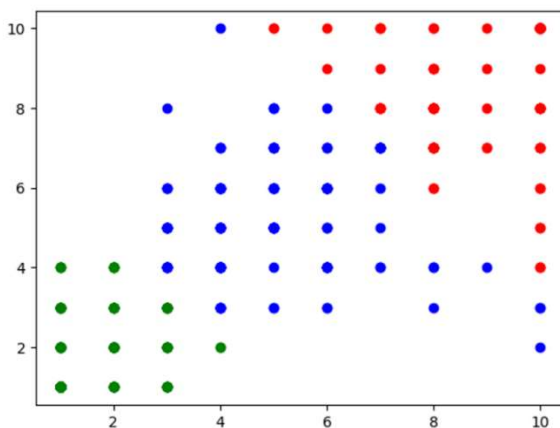
$k = 3$: ECR = 6.67

b) $k = 2$: silhouette = 0.5968

$k = 3$: silhouette = 0.5245

A silhueta é um critério interno de medição que avalia a separação e a coesão da solução de clustering, enquanto o ECR é um critério externo que avalia o quão bem os clusters se ajustam à classificação real das observações. Assim, comparando as silhuetas verifica-se que com $k = 2$ obtêm-se clusters melhor separados e mais coesos, mas, tendo em conta o ECR, verificamos que estes agrupamentos não são homogêneos, sendo, de acordo com esta métrica, a solução com $k = 3$ melhor.

5)



6) Ao observarmos o gráfico obtido, verificamos que as 2 features com maior mutual information produzem uma boa solução de clustering, uma vez que se observa uma boa coesão dos pontos pertencentes ao mesmo cluster e uma boa separação entre os clusters.

III. APPENDIX

```
import numpy as np
import scipy.io.arff
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
from sklearn.feature_selection import SelectKBest, mutual_info_classif

arff = scipy.io.arff

breast_data, breast_meta = arff.loadarff("breast.w.new.arff")

breast_targets = np.array(list(el["Class"] for el in breast_data))
breast_X = np.array(np.array(list(list(el[i] for i in range(9)) for el in breast_data)))

clusters2 = KMeans(n_clusters=2)
clusters3 = KMeans(n_clusters=3)

clusters2.fit(breast_X)
clusters3.fit(breast_X)

#4a

benign2 = [0, 0]
benign3 = [0, 0, 0]
malign2 = [0, 0]
malign3 = [0, 0, 0]

for i in range(len(clusters2.labels_)):
    if breast_targets[i] == b'benign':
        benign2[clusters2.labels_[i]] += 1
    else:
        malign2[clusters2.labels_[i]] += 1

for i in range(len(clusters3.labels_)):
    if breast_targets[i] == b'benign':
        benign3[clusters3.labels_[i]] += 1
    else:
        malign3[clusters3.labels_[i]] += 1

total_cluster_2_0 = benign2[0] + malign2[0]
total_cluster_2_1 = benign2[1] + malign2[1]

total_cluster_3_0 = benign3[0] + malign3[0]
total_cluster_3_1 = benign3[1] + malign3[1]
total_cluster_3_2 = benign3[2] + malign3[2]
```

Aprendizagem 2021/22
Homework IV – Group 17

```
ecr2 = 0.5 * ((total_cluster_2_0 - max(benign2[0], malign2[0])) + (total_cluster_2_1 -  
max(benign2[1], malign2[1])))  
ecr3 = (1/3) * ((total_cluster_3_0 - max(benign3[0], malign3[0])) + (total_cluster_3_1 -  
max(benign3[1], malign3[1])) + (total_cluster_3_2 - max(benign3[2], malign3[2])))  
  
print(ecr2, ecr3)  
  
#4b  
  
silhouette2 = silhouette_score(breast_X, clusters2.labels_)  
silhouette3 = silhouette_score(breast_X, clusters3.labels_)  
  
print(silhouette2, silhouette3)  
  
#5  
  
kbest = SelectKBest(mutual_info_classif, k=2).fit_transform(breast_X, breast_targets)  
  
clusters = KMeans(n_clusters=3)  
clusters.fit(kbest)  
  
clusters_div = [[], [], []]  
  
for i in range(len(clusters.labels_)):  
    clusters_div[clusters.labels_[i]].append(kbest[i])  
  
for i in range(len(clusters_div)):  
    clusters_div[i] = np.array(clusters_div[i])  
  
plt.scatter(clusters_div[0][:,0], clusters_div[0][:,1], color="red")  
plt.scatter(clusters_div[1][:,0], clusters_div[1][:,1], color="green")  
plt.scatter(clusters_div[2][:,0], clusters_div[2][:,1], color="blue")  
  
plt.savefig("graphic.png", format="png")
```

END