

Francisco Silva, 110409
Joana Lardoso, 109864

Homework 1

1. Pen-and-paper

$$1) IG(\text{class}, Y_i) = H(\text{class}) - H(\text{class}|Y_i)$$

Para $Y_1 > 0,4$

$$H(\text{class}) = - \sum_{x \in \text{class}} P(\text{class} = x) \cdot \log_2(P(\text{class} = x)) =$$

$$= -(0,25 \times \log_2(0,25) + (0,375 \times \log_2(0,375))_2) = 1,562 \text{ bits}$$

$$H(\text{class}|Y_2) = \frac{1}{8}(-(0+1\log 1+0)) + \frac{4}{8}(-(\frac{1}{2}\log \frac{1}{2} + (\frac{1}{4}\log \frac{1}{4}) + 2)) + \\ + \frac{3}{8}(-(0+2(\frac{1}{3}\log \frac{1}{3}) + (\frac{2}{3}\log \frac{2}{3}))) =$$

$$= 1,094$$

$$H(\text{class}|Y_3) = -\frac{3}{8} \times (3(\frac{1}{3}\log \frac{1}{3})) - \frac{2}{8} \times (2(\frac{1}{2}\log \frac{1}{2})) -$$

$$-\frac{3}{8}(\frac{3}{4}\log \frac{3}{4} + \frac{1}{4}\log \frac{1}{4}) = 1,439$$

$$H(\text{class}|Y_4) = -\frac{1}{8}(1\log 1) - \frac{3}{8}(\frac{2}{3}\log \frac{2}{3} + \frac{1}{3}\log \frac{1}{3}) - \frac{4}{8}(0 + \frac{1}{4}\log \frac{1}{4}) + \frac{3}{4}\log \frac{3}{4}$$

$$= \cancel{1,094} 0,75$$

$$IG(\text{class}, Y_2) = H(\text{class}) - H(\text{class}|Y_2) = 0,468$$

$$IG(\text{class}, Y_3) = H(\text{class}) - H(\text{class}|Y_3) = 0,123$$

$$IG(\text{class}, Y_4) = H(\text{class}) - H(\text{class}|Y_4) = 0,812$$

We choose Y_4 to be the next node because it has the highest Information Gain

Para: $Y_1 \leq 0,4 \wedge Y_4 = 0$:

mode(~~A, B~~(B)) = B (alphabetic order)

Para: $Y_1 \leq 0,4 \wedge Y_4 = 1$:

mode(A, A, B) = A

Para: $Y_1 \geq 0,4 \wedge Y_4 = 2$: Has ~~more than~~ four (4) observations,
therefore it will have to be expanded

| Para $Y_1 \geq 0,4 \wedge Y_4 = 2$

$$H(\text{class}) = -\left(\frac{1}{4} \log \frac{1}{4} + \frac{3}{4} \log \frac{3}{4}\right) = 0.811$$

$$H(\text{class}|Y_2) = -\frac{1}{4}(1 \log 1) - \frac{3}{4}\left(\frac{1}{3} \log \frac{1}{3} + \frac{2}{3} \log \frac{2}{3}\right) = 0.689$$

$$H(\text{class}|Y_3) = -\frac{1}{4}(1 \times \log 1) - \frac{2}{4}\left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2}\right) = 0.5$$

$$IG(\text{class}, Y_2) = 0.811 - 0.689 = 0.121$$

$$IG(\text{class}, Y_3) = 0.811 - 0.5 = 0.311 \leftarrow$$

We choose Y_3 to be the next node because it has the highest information gain

Para $Y_1 \geq 0,4 \wedge Y_4 = 2 \wedge Y_3 = 0$

mode(C) = C

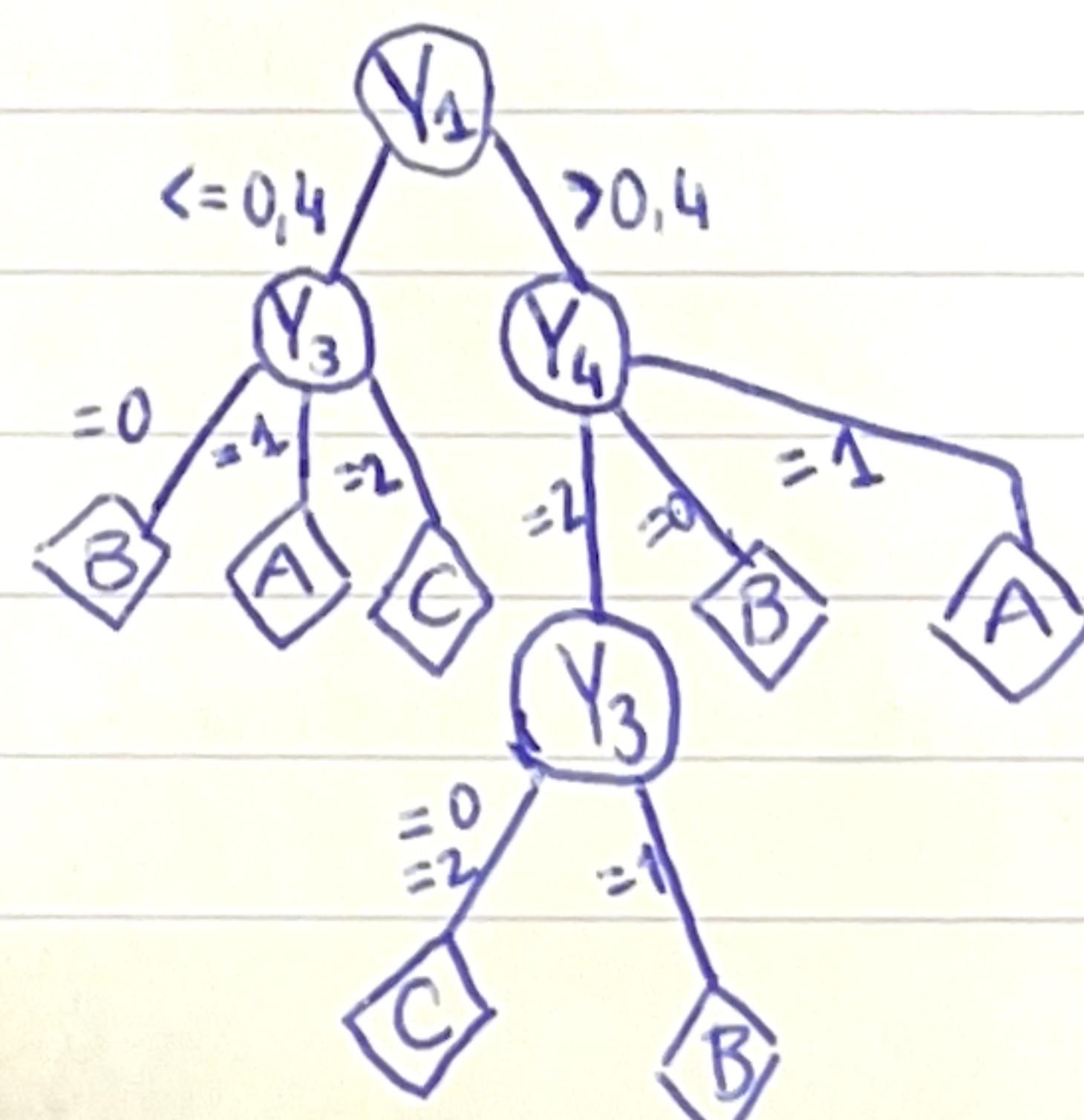
Para $Y_1 \geq 0,4 \wedge Y_4 = 2 \wedge Y_3 = 1$

mode(C, B) = B (alphabetic order)

Para $Y_1 \geq 0,4 \wedge Y_4 = 2 \wedge Y_3 = 2$

mode(C) = C

Assim, obtenemos a arvore:



2) Confusion Matrix

		True		
		A	B	C
Predictions	A	4	1	0
	B	0	3	1
	C	0	0	3

$$\rightarrow \begin{bmatrix} 4 & 1 & 0 \\ 0 & 3 & 1 \\ 0 & 0 & 3 \end{bmatrix}$$

3) TP: True Positive (Previsões corretas) real = classe; Previsão = classe

FP: False Positive (Previsões erradas) real ≠ classe; Previsão = classe

TN: True Negative (Falha em prever) real = classe; Previsão ≠ classe

FN: False Negative (Previsões corretas de não-classe) real ≠ classe; previsão ≠ classe

$$F_1 = \frac{2 \times P \times R}{P+R}$$

$$P = \underbrace{\frac{TP}{TP+FP}}_{\text{Precision}}$$

$$R = \underbrace{\frac{TP}{TP+FN}}_{\text{Recall / sensitivity}}$$

Classe A : $TP_A = 4$ $FP_A = 1$ $FN_A = 0$

$$P_A = \frac{4}{4+1} = \frac{4}{5} \quad R_A = \frac{4}{4+0} = 1 \quad F_1 = \frac{2 \times 0.8 \times 1}{1+0.8} = 0.89$$

Classe B : $TP_B = 3$ $FP_B = 1$ $FN_B = 1$

$$P_B = \frac{3}{3+1} = \frac{3}{4} \quad R_B = \frac{3}{3+1} = \frac{3}{4} \quad F_1 = \frac{2 \times 0.75 \times 0.75}{1.5} = \frac{3}{4} = 0.75$$

Classe C :

$$P_C = \frac{3}{3} = 1 \quad R_C = \frac{3}{3+1} = \frac{3}{4} \quad F_1 = \frac{2 \times \frac{3}{4}}{1.75} = 0.86$$

Class B is the one with the lowest training score

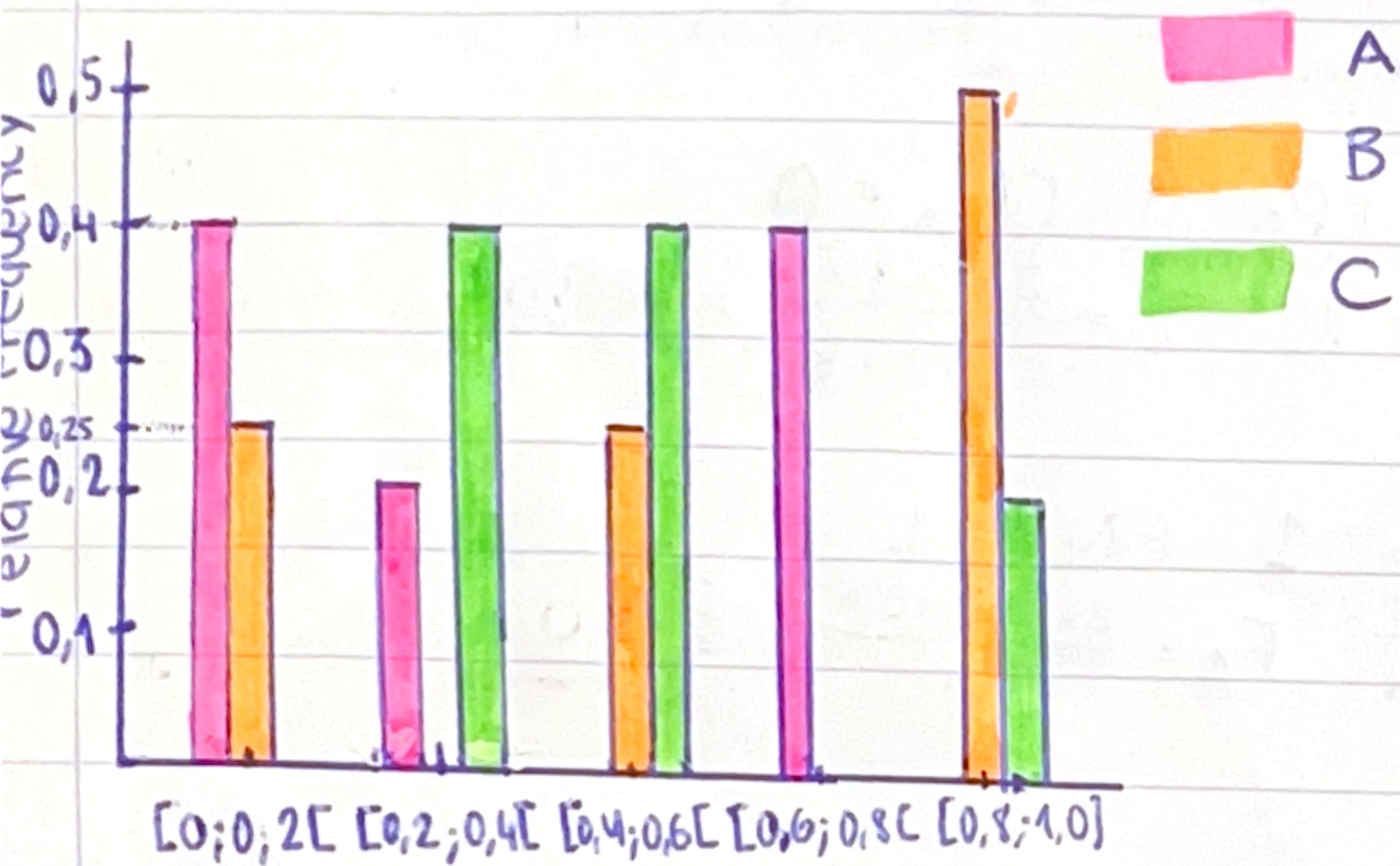
4)

- A: 0,12 0,33 0,62 0,71
 B: 0,18 0,52 0,83 0,90
 C: 0,25 0,45 0,58 0,95

$[0,00; 0,20]$ | $[0,20; 0,40]$ | $[0,40; 0,60]$ | $[0,60; 0,80]$ | $[0,80; 1]$
 AAB ACC BCC AA BBC

- A: $[0,40; 0,20; 0; 0,40; 0]$
 B: $[0,25; 0; 0,25; 0; 0,5]$
 C: $[0; 0,40; 0,40; 0; 0,20]$

Histograma A B e C

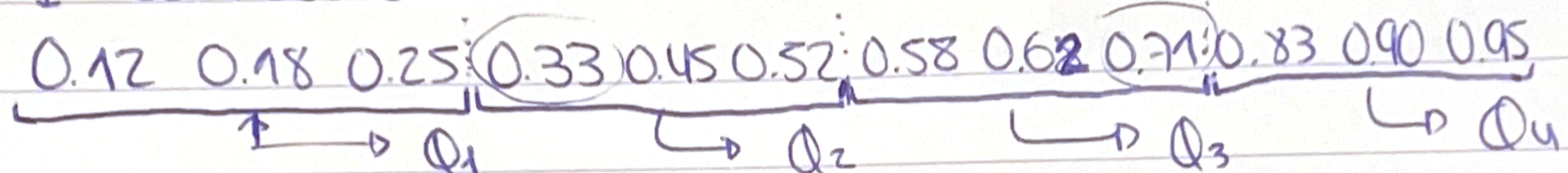


domina: $\left\{ \begin{array}{l} A \text{ em } [0,0; 0,20] \wedge [0,60, 0,80] \\ C \text{ em } [0,20; 0,60] \\ B \text{ em } [0,80; 1] \end{array} \right.$

Assim, n é 4 (dividimos o intervalo $[0,1]$ em 4)

5. Neste dataset D, para encontrar outliers temos de utilizar o método Interquartile Range (IQR), aplicado apenas a variáveis [continuas, logo V₁. numérica]

1º Sort dos dados



$$\text{IQR} = \max(Q_3) - \min(Q_1) = 0.71 - 0.33 \\ = 0.38$$

intervalo dos valores que não são outliers: $[Q_1 - 1,5 \times \text{IQR}, Q_3 + 1,5 \times \text{IQR}] =$

$$= \left[\frac{0.25 + 0.33}{2} - 1,5 \times 0.38, \frac{0.71 + 0.83}{2} + 1,5 \times 0.38 \right] = [-0.28; 1.34]$$

Como todos os dados pertencem ao intervalo [-0.28; 1.34] podemos concluir que no dataset D não há outlier, pelo IQR method.