



homework

# Homework 1:

## I. Pen-and-Paper:

1)

$$\text{Shannon Entropy : } I(n) = \sum_{n \in X} p(n) \times I(n) = - \sum_{n \in X} p(n) \times \log_2(p(n))$$

$$E(\text{class}) = E(y_{\text{out}}) = - \left( \frac{1}{3} \log_2 \left( \frac{1}{3} \right) + \frac{1}{3} \log_2 \left( \frac{1}{3} \right) + \frac{1}{3} \log_2 \left( \frac{1}{3} \right) \right) \\ = \log_2 3 = 1.58496$$

$$IG(\text{class} | y_n) = E(\text{class}) - E(\text{class} | y_n)$$

$$y_1 > 0.4$$

$$E(\text{class} | y_2) = E(y_{\text{out}} | y_2=0) + E(y_{\text{out}} | y_2=1) + E(y_{\text{out}} | y_2=2) \\ = \frac{3}{7} \left( - \left( \frac{1}{3} \log_2 \left( \frac{1}{3} \right) + \frac{1}{3} \log_2 \left( \frac{1}{3} \right) + \frac{1}{3} \log_2 \left( \frac{1}{3} \right) \right) \right) + \\ + \frac{2}{7} \left( - \left( \frac{1}{2} \log_2 \left( \frac{1}{2} \right) + \frac{1}{2} \log_2 \left( \frac{1}{2} \right) \right) + \frac{2}{7} \left( - \left( 1 \log_2 (1) \right) \right) \right) \\ = \frac{3}{7} \left( - \log_2 \left( \frac{1}{3} \right) \right) + \frac{2}{7} \left( - \log_2 \left( \frac{1}{2} \right) \right) + \frac{2}{7} \left( - \log_2 (1) \right) \\ = \frac{3}{7} \log_2 (3) + \frac{2}{7} \approx 0.964984$$

$$E(\text{class} | y_3) = E(y_{\text{out}} | y_3)$$

$$= E(y_{\text{out}} | y_3=0) + E(y_{\text{out}} | y_3=1) + E(y_{\text{out}} | y_3=2) \\ = \frac{1}{7} \left( - 1 \log_2 (1) \right) + \frac{2}{7} \left( - \left( \frac{1}{2} \log_2 \left( \frac{1}{2} \right) + \frac{1}{2} \log_2 \left( \frac{1}{2} \right) \right) \right) \\ + \frac{4}{7} \left( - \left( \frac{1}{2} \log_2 \left( \frac{1}{2} \right) + \frac{1}{2} \log_2 \left( \frac{1}{2} \right) \right) \right) \\ = \frac{1}{7} \log_2 (1) + \frac{2}{7} + \frac{4}{7} = \frac{6}{7} \approx 0.857143$$

$$E(\text{Class1} | y_4) = E(y_{\text{out}} | y_4) = E(y_{\text{out}} | y_4=0) + E(y_{\text{out}} | y_4=1) + E(y_{\text{out}} | y_4=2)$$

$$\begin{aligned}
 &= \frac{2}{7} \left( -\left( \frac{1}{2} \log_2 \left( \frac{1}{2} \right) + \frac{1}{2} \log_2 \left( \frac{1}{2} \right) \right) \right) + \\
 &+ \frac{3}{7} \left( -\left( \frac{1}{3} \log_2 \left( \frac{1}{3} \right) + \frac{2}{3} \log_2 \left( \frac{2}{3} \right) \right) + \frac{2}{7} \left( -\left( \frac{1}{2} \log_2 \left( \frac{1}{2} \right) + \frac{1}{2} \log_2 \left( \frac{1}{2} \right) \right) \right) \right) \\
 &= \frac{2}{7}(1) + \frac{3}{7} \left( -\left( \frac{1}{3} (2 - 3) \log_2 (3) \right) \right) + \frac{2}{7} \\
 &= \frac{2}{7} + \frac{3}{7} \log_2 (3) \approx 0.964984
 \end{aligned}$$

$$IG(\text{Class1} | y_2) = E(\text{Class1}) - E(\text{Class1} | y_2)$$

$$= 1.58496 - 0.964984 = 0.619976$$

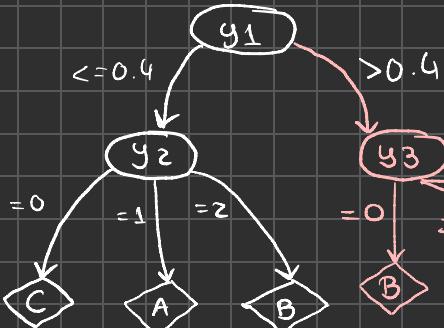
$$IG(\text{Class1} | y_3) = E(\text{Class1}) - E(\text{Class1} | y_3)$$

$$= 1.58496 - 0.857143 = 0.727817$$

$$IG(\text{Class1} | y_4) = E(\text{Class1}) - E(\text{Class1} | y_4)$$

maior IG, escolhe esta variável ( $y_3$ )

$$= 1.58496 - 0.964984 = 0.619976$$



- Como temos 4 observações para  $y_1 > 0.4 \wedge y_3 = 2$ , temos de recorrer novamente a IG.

- Para  $y_1 > 0.4$  e  $y_3 = 0$ , há apenas 1 observação, pelo que o valor da variável de destino será B. Nesta situação (única observação)

- Para  $y_1 > 0.4$  e  $y_3 = 1$ , há apenas 2 observações, pelo que o NC não é dividido. Segundo a ordem alfabética, o valor da variável de destino será A.

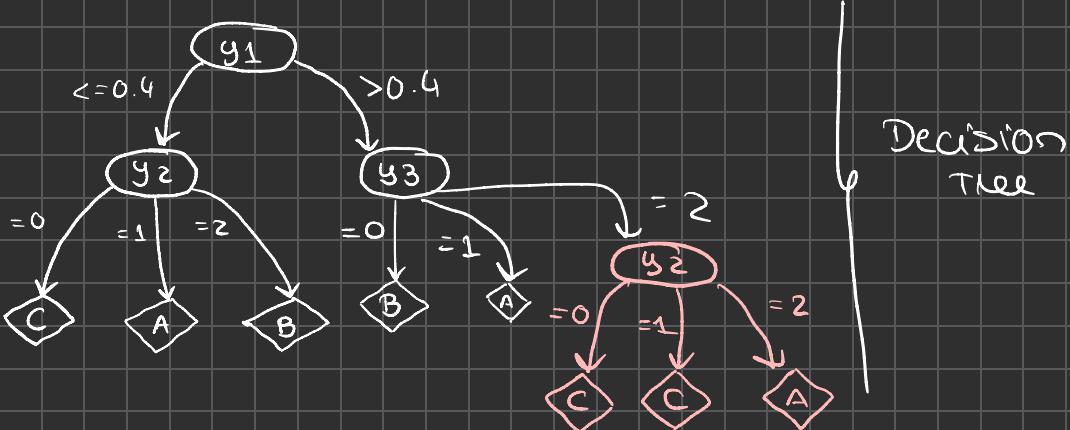
$$\begin{aligned}
 E(\text{class} | y_3 = 2 \wedge y_2) &= E(\text{class} | y_3 = 2 \wedge y_2 = 0) + E(\text{class} | y_3 = 2 \wedge y_2 = 1) + \\
 &\quad + E(\text{class} | y_3 = 2 \wedge y_2 = 2) \\
 &= \frac{1}{4} \left( -\underbrace{1 \log_2(1)}_{\emptyset} \right) + \frac{1}{4} \left( -\underbrace{1 \log_2(1)}_{\emptyset} \right) + \frac{1}{2} \left( -\underbrace{1 \log_2(1)}_{\emptyset} \right) \\
 &= \boxed{0}
 \end{aligned}$$

$$\begin{aligned}
 E(\text{class} | y_3 = 2 \wedge y_4) &= E(\text{class} | y_3 = 2 \wedge y_4 = 0) + E(\text{class} | y_3 = 2 \wedge y_4 = 1) + \\
 &\quad + E(\text{class} | y_3 = 2 \wedge y_4 = 2) \\
 &= \frac{1}{2} \left( -\left( \frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right) + \right. \\
 &\quad \left. + \frac{1}{4} \left( -\underbrace{1 \log_2(1)}_{\emptyset} \right) + \frac{1}{4} \left( -\underbrace{1 \log_2(1)}_{\emptyset} \right) \right) \\
 &= \frac{1}{2} (-(-1)) = \frac{1}{2} = \boxed{0.5}
 \end{aligned}$$

$$\begin{aligned}
 IG(\text{class} | y_3 = 2 \wedge y_2) &= E(\text{class}) - \underbrace{E(\text{class} | y_3 = 2 \wedge y_2)}_{\emptyset} \\
 &= E(\text{class}) = \boxed{1.58496}
 \end{aligned}$$

*G maior IG, excolher-se  
esta variável*

$$\begin{aligned}
 IG(\text{class} | y_3 = 2 \wedge y_4) &= E(\text{class}) - E(\text{class} | y_3 = 2 \wedge y_4) \\
 &= 1.58496 - 0.5 \\
 &= \boxed{1.08496}
 \end{aligned}$$



Para  $y_1 > 0.4 \wedge y_3 = 2$  já só temos 4 observações: 1 observação para  $y_2 = 0$ , 1 observação para  $y_2 = 1$  e 2 observações para  $y_2 = 2$ . Desta forma os nós não serão expandidos. Para os ramos  $y_2 = 0$  e  $y_2 = 1$ , como cada um tem apenas uma observação, o valor da variável de destino de cada ramo será o valor da observação (para ambos, C). Para o ramo  $y_2 = 2$ , ambas as observações têm valor A, pelo que a variável de destino deste ramo tem valor A.

2.

		predicto		
		A	B	C
real	A	4	0	0
	B	1	2	1
	C	0	0	4

		real		
		A	B	C
predicto	A	4	1	0
	B	0	2	0
	C	0	1	4

3.

which class has the lower F1-score?

$$F1\text{-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Sensitivity / Recall} = \frac{TP}{TP + FN}$$

Class A:

$$\text{Precision}_A = \frac{4}{4+0+1} = \frac{4}{5} = 0.80 \quad \text{Recall}_A = \frac{4}{4} = 1$$

$$F1\text{-score}_A = 2 \times \frac{0.80 \times 1}{0.80 + 1} = 2 \times \frac{0.80}{1.80} = \frac{1.60}{1.80} = \frac{8}{9} = 0.8(8)$$

Class B:

$$\text{Precision}_B = \frac{2}{2+0} = 1 \quad \text{Recall}_B = \frac{2}{4} = \frac{1}{2} = 0.5$$

$$F1\text{-score}_B = 2 \times \frac{1 \times 0.5}{1+0.5} = 2 \times \frac{0.5}{1.5} = 2 \times \frac{1}{3} = \frac{2}{3} = 0.6(6)$$

TP - true positive

FP - false positive

TN - true negative

FN - false negative

## class C :

$$\text{Precision}_C = \frac{4}{4+1+0} = \frac{4}{5} = 0.8 \quad \text{Recall}_C = \frac{4}{4} = 1$$

$$F1 - Score_C = 2 \times \frac{0.8 \times 1}{0.8 + 1} = \frac{1.60}{1.80} = \frac{8}{9} = 0.8(8)$$

R: Class B is the one with the lowest training F1 score, with a value of approximately 0.67.

4.

	rank y <sub>1</sub>	rank y <sub>2</sub>	$(\text{rank } y_1)^2$	$(\text{rank } y_2)^2$	rank y <sub>1</sub> · rank y <sub>2</sub>
x <sub>1</sub>	3	8	9	64	24
x <sub>2</sub>	2	11	4	121	22
x <sub>3</sub>	1	3,5	1	12,25	3,5
x <sub>4</sub>	5	3,5	25	12,25	17,5
x <sub>5</sub>	4	3,5	16	12,25	14
x <sub>6</sub>	10	11	100	121	110
x <sub>7</sub>	12	3,5	144	12,25	42
x <sub>8</sub>	11	11	121	121	121
x <sub>9</sub>	7	8	49	64	56
x <sub>10</sub>	9	3,5	81	12,25	31,5
x <sub>11</sub>	6	8	36	64	48
x <sub>12</sub>	8	3,5	64	12,25	28
$\sum_{x \in X} (\text{rank } y_1) =$	$\sum_{x \in X} (\text{rank } y_2) =$	$\sum_{x \in X} (\text{rank } y_1)^2 =$	$\sum_{x \in X} (\text{rank } y_2)^2 =$	$\sum_{x \in X} (\text{rank } y_1)(\text{rank } y_2) =$	
= 78	= 78	= 650	= 628,5	= 517,5	

Há rank ties, pelo que o clássico Spearman é substituído pelo PCC dos ranks.

$$\text{Spearman}(y_1, y_2) = \text{PCC}((3, 2, 1, 5, 4, 10, 12, 11, 7, 9, 6, 8), \\ (8, 11, 3.5, 3.5, 3.5, 11, 3.5, 11, 8, 3.5, 8, 3.5))$$

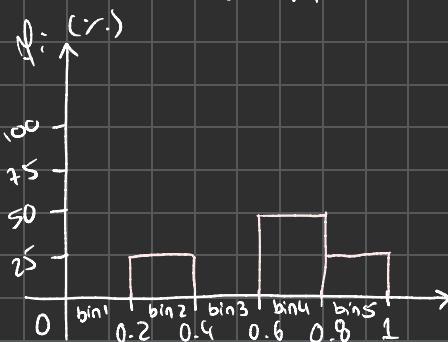
$$n = \frac{\sum(\text{rank } y_1)(\text{rank } y_2) - \frac{\sum(\text{rank } y_1)}{n} \sum(\text{rank } y_2)}{\sqrt{\left( \sum(\text{rank } y_1)^2 - \frac{(\sum(\text{rank } y_1))^2}{n} \right) \cdot \left( \sum(\text{rank } y_2)^2 - \frac{(\sum(\text{rank } y_2))^2}{n} \right)}}$$

$$r = \frac{517.5 - \frac{78 \times 78}{12}}{\sqrt{\left(650 - \frac{78^2}{12}\right)\left(628.5 - \frac{78^2}{12}\right)}} \approx 0.0796587$$

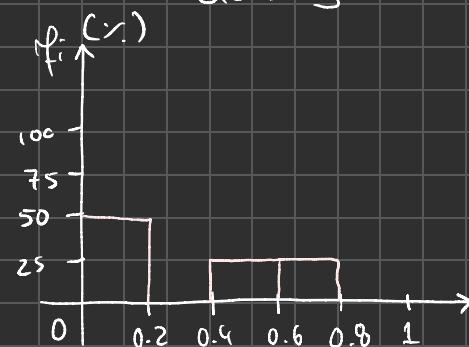
R: Como  $r = 0.0796587$ ,  $y_1$  e  $y_2$  são ligeiramente relacionados.

5.

classe A :



classe B :



classe C :

