

LEIC-T 2023/2024
Aprendizagem - Machine Learning
Homework I
Deadline 29/9/2024 20:00
Submit on Fenix as pdf

I) Correlation (2 pts)

Compute the correlation (Pearson correlation) and **Spearman's rank** for two variables x_1 and x_2

(a) (1 pts)

$x_1 = (1, 3, 4, 6)$

$x_2 = (-30, -10, 0, 20)$

Solution, we assume it. Is a sample

(if you assume it is a population you will get the same results for correlation):

$Mean(x_1)=7/2, Mean(x_2)=-5$

Covariance for Sample $cov(x_1, x_2)=130/3$

Standard Deviation for sample: $s_1=Sqrt(13/3), s_2=10 * Sqrt(13/3)$

$corr(x_1, x_2)=cov(x_1, x_2)/(s_1*s_2)=1$

$R_1=(1, 2, 3, 4)$

$R_2=(1, 2, 3, 4)$

$Mean(R_1)=5/2, Mean(R_2)=5/2$

Covariance for Sample $cov(R_1, R_2)=5/3$

Standard Deviation for sample: $s_1=Sqrt(5/3), s_2=Sqrt(5/3)$

$corr(R_1, R_2)=cov(R_1, R_2)/(s_1*s_2)=1$

Why are the same? *Points are on the line.*

(b) (1pts)

$x_1 = (1, 3, 4, 6)$

$x_2 = (-3, -0.5, 29, 30)$

$Mean(R_1)=7/2, Mean(R_2)=13.87$

Covariance for Sample $cov(R_1, R_2)=32.4167$

Standard Deviation for sample: $s_1=Sqrt(13/3), s_2=18.075$

$corr(x_1, x_2)=cov(x_1, x_2)/(s_1*s_2)=0.861516$

LEIC-T 2023/2024
 Aprendizagem - Machine Learning
 Homework I
 Deadline 29/9/2024 20:00
Submit on Fenix as pdf

$R1=(1,2,3,4)$

$R2=(1,2,3,4)$

Rank is the same as before...

$$\text{corr}(R1,R2)=\text{cov}(R1,R2)/(s1*s2)=1$$

Why are they different? *Points are ordered but are not on the line.*

II) Decision Trees (5 pts)

$F1$	$F2$	$F3$	$F4$	$Output$
c	a	b	x	n
a	a	c	a	t
a	b	b	a	t
c	b	c	x	m
a	b	b	c	f

(a) (2 pts)

Determine the root of decision tree using the ID3 algorithm with the target “Output”. Indicate the calculation.

Solution:

$$p(n)=1/5 \quad p(t)=2/5, \quad p(m)=1/5, \quad p(f)=1/5$$

$$\text{Log}_2[x] = \text{Log}[x]/\text{Log}[2]$$

$$I(\text{table}) = -3*1/5*\text{Log}_2[1/5] - 2/5*\text{Log}_2[2/5] = 1.92193 \text{ bits}$$

$$E(P) = \sum_{i=1}^n \frac{|C_i|}{|C|} I(C_i) \quad \text{gain}(P) = I(C) - E(P)$$

$$F1 \quad Ca=(t,t,f), \quad Cc=Cx=(n,m),$$

$$I(Ca) = -2/3*\text{Log}_2[2/3] - 1/3*\text{Log}_2[1/3] = 0.918296 \text{ bit}$$

$$I(Cc) = I(Cx) = -1/2*\text{Log}_2[1/2] - 1/2*\text{Log}_2[1/2] = 1 \text{ bit}$$

$$E(F1) = 2/5*1 + 3/5*0.918296 = 0.950978 \text{ bit}$$

$$\text{Gain}(F1) = \text{Gain}(F4) = 1.92193 - 0.950978 = 0.970952 \text{ bit}$$

$$F2 \quad Ca=(n,t) \quad Cb=(t,m,f)$$

$$I(Ca) = -1/2*\text{Log}_2[1/2] - 1/2*\text{Log}_2[1/2] = 1 \text{ bit}$$

$$I(Cb) = -3*1/3*\text{Log}_2[1/3] = 1.58496 \text{ bit}$$

LEIC-T 2023/2024
 Aprendizagem - Machine Learning
 Homework I
 Deadline 29/9/2024 20:00

Submit on Fenix as pdf

$$E(F2) = 2/5 * 1 + 3/5 * 1.58496 = 1.35098 \text{ bit}$$

$$\text{Gain}(F2) = 1.92193 - 1.35098 = 0.57095 \text{ bit}$$

$$F3 \text{ } Cb=(n,t,f) \text{ } Cc=(t,m)$$

$$I(Cb) = -3 * 1/3 * \log_2[1/3] = \log_2[3] = 1.58496 \text{ bit}$$

$$I(Cc) = -2 * 1/2 \log_2[1/2] = \log_2[2] = 1 \text{ bit}$$

$$E(F3) = 3/5 * 1.58496 + 2/5 * 1 = 1.3509 \text{ bit}$$

$$\text{Gain}(F3) = 1.92193 - 1.3509 = 0.57103 \text{ bit}$$

$$F4 \text{ } Ca=(t,t) \text{ } Cx=(n,m) \text{ } Cc=(f)$$

$$I(Ca) = I(Cc) = 0$$

$$I(Cx) = -2 * 1/2 \log_2[1/2] = \log_2[2] = 1 \text{ bit}$$

$$E(F4) = 2/5 * 1 + 2/5 * 0 + 1/5 * 0 = 0.4 \text{ bit}$$

$$\text{Gain}(F4) = 1.92193 - 0.4 = 1.52193 \text{ bit}$$

*We chose **F4** as the root*

(b) (2 pts)

Determine the decision tree using the ID3 algorithm with the target “Output”. Indicate the calculation and draw your decision tree.

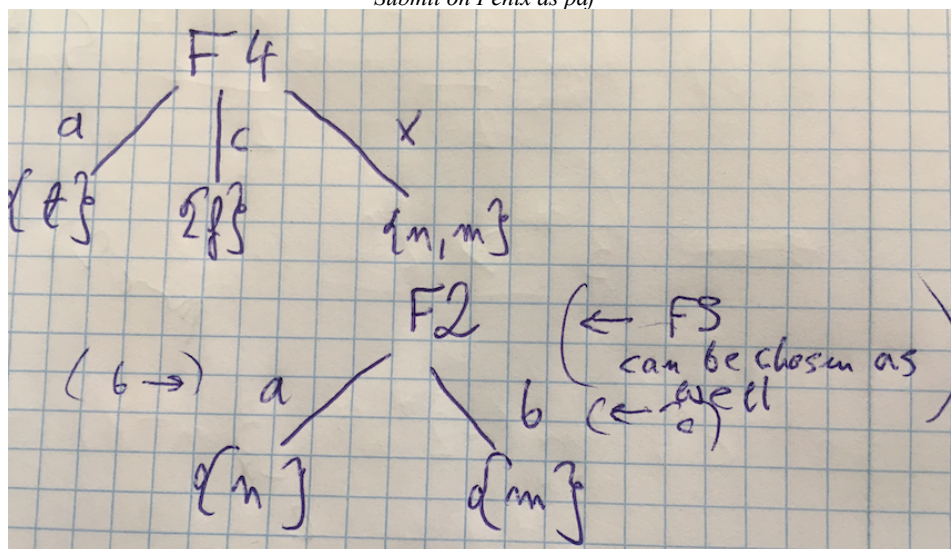
Solution:

The remaining table Cx:

<u>F1</u>	<u>F2</u>	<u>F3</u>	<u>Output</u>
c	a	b	n
c	b	c	m

F2 and F3 are equal, they give the same gain. We chose for the tree F2

LEIC-T 2023/2024
Aprendizagem - Machine Learning
Homework I
Deadline 29/9/2024 20:00
Submit on Fenix as pdf



(c) (1 pts)

Draw the training confusion matrix for the learnt decision tree.

	TRUE				
		n	t	m	f
PREDICTED	n	1	0	0	0
	t	0	2	0	0
	m	0	0	1	0
	f	0	0	0	1