

Data Mining 2020 Final

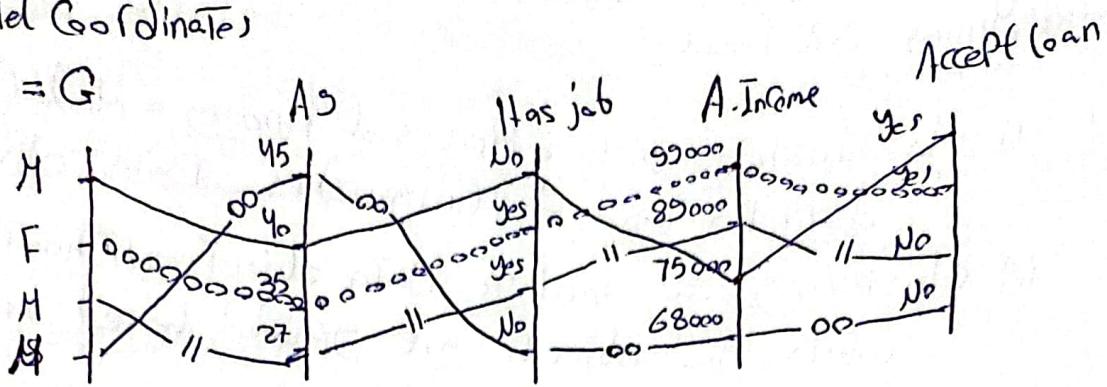
- Q1) i) Classification \Rightarrow is a process of finding a Model That Describes and distinguishes Data classes
- ii) Clustering \Rightarrow analyze Data objects without class labels, The objects are grouped based on similarity
- iii) Dendrogram \Rightarrow Graph That Show us How we can Merge Two object Together and we can construct The resulted clusters after doing Agglomerative Algo
- iv) Classification Vs Clustering
- Dividing Data Set based on Values of Class labels
 - Predictive Supervised learning
 - Dividing Data set based on Similarity
 - Descriptive
 - UnSupervised learning
- v) Model Construction \Rightarrow is first process of classification That use Training Data set To construct Model, This Model help us To predict The value of class label of New object
- vi) Model evaluation \Rightarrow is second process of classification That Determiner How The Model is good for predict The CL Value ~~of test~~ on The ~~Set of~~ Set of records of Testing Data set
- vii) Model usage \Rightarrow use Model To predict a new CL Value of New object
- viii) Evaluation Matrix \Rightarrow This Matrix help w To Compute accuracy & precision & Recall of classifier To determine How This classifier is good

| | |
|----|----|
| TP | FN |
| FP | TN |

11

b) i) Parallel Coordinates

$$G_{MLR} = G$$



ii) $(obj_1 \text{ } \& \text{ } obj_4)$

$$obj_1 = (M, 40, No, 75000, Yes)$$

$$obj_4 = (M, 45, No, 68000, No)$$

$$\delta_G = 1, \delta_{Age} = 1, \delta_{Has\ job} = 1, \delta_{Income} = 1, \delta_{Accept\ loan} = 1$$

$$Dis_G = 0 \quad (M=M)$$

$$Dis_{Age} = \frac{140 - 45}{45 - 27} = 0.27$$

$$Dis_{Income} = \frac{175000 - 68000}{99000 - 68000} = 0.22$$

$$Dis_{Has\ job} = 0 \quad (No=No)$$

$$Dis_{Accept\ loan} = 1$$

$$Dis(obj_1, obj_2) = \frac{1 \times 0 + 1 \times 0.27 + 1 \times 0 + 1 \times 0.22 + 1 \times 1}{1+1+1+1+1} = 0.298$$

$$Sim(obj_1, obj_2) = 1 - Dis = 1 - 0.298 = 0.702$$

$$obj_2 = (F, 35, Yes, 99000, Yes)$$

$$obj_3 = (M, 27, Yes, 89000, No)$$

$(obj_2 \text{ } \& \text{ } obj_3)$

$$Dis_G = 1 \quad (F \neq M)$$

$$Dis_{Age} = \frac{135 - 27}{45 - 27} = \frac{1}{9}$$

$$Dis_{Income} = \frac{199000 - 89000}{99000 - 68000} = \frac{10}{31}$$

$$Dis_{Has\ job} = 0$$

$$Dis_{Accept\ loan} = 1$$

$$Dis(obj_2, obj_3) = \frac{1 \times 1 + 1 \times \frac{1}{9} + 1 \times 0 + 1 \times \frac{10}{31} + 1 \times 1}{1+1+1+1+1} = 0.553$$

$$Sim(obj_2, obj_3) = 1 - Dis = 1 - 0.553 = 0.447$$

②

Q2

$$= \frac{25+27+32+45+43+28+29+17+21+31+26+42}{12} = 31,3 \approx 31$$

12 pos 31
and also Credit Rating is (Missing) loop

Misstry \Rightarrow Fair ^($\text{bad} \rightarrow \text{good}$) fair = 8
~~excellent~~ excellent = 15

Student is Missing Notify

Missing = Yes No = 6
Yes = 7

No = 6
Yes = 7

b) New object (20, low, yes, fair, ?) ~~test~~

b) New object (0, low, yes, ...),
b) New object (label)'s No, Yes (With 2 class labels)
و طبقاً (نعم و لا) + انفع هو

$$P(\text{yes}) = \frac{9}{14}$$

$$P(\text{No}) = \frac{5}{14}$$

جديد كإبلا للـ New object اجب لا Yes من و No من

$$P(20 \mid \text{Yes})$$

١٥
Cor-
ATII
العاصمة
والمجتمع

$$\mu_{\text{age}} = \frac{32 + 31 + 45 + 31 + 29 + 17 + 21 + 31 + 36}{9} = 30.33$$

$$G(\text{age}) = \sqrt{\frac{(32 - 30.3)^2 + (31 - 30.3)^2 + \dots + (36 - 30.3)^2}{9}} = 7.58$$

$$P(20 \text{ Yes}) = \frac{1}{7.58 * \sqrt{2\pi}} * e^{-\frac{20^2}{7.58^2}} = 0.020$$

$$P(\text{low}|\text{yes}) = \frac{3}{9} = \frac{1}{3}$$

$$P(\text{yes}|\text{yes}) = \frac{7}{9}$$

$$P(\text{fair}|\text{yes}) = \frac{6}{9} = \frac{2}{3}$$

$$P(\text{yes|New}) = P(\text{?o|yes}) * P(\text{low|yes}) * P(\text{yes|yes}) * P(\text{fair|yes})$$

$$= 0.020 * \frac{1}{3} * \frac{7}{9} * \frac{2}{3} * \frac{6}{9} = 2.32 * 10^{-3}$$

Yes \rightarrow New \rightarrow New \rightarrow Probabil
No \rightarrow New \rightarrow No \rightarrow Probabil

P(?o|No)

$$\mu_{(\text{ave})} = \frac{25+27+43+28+42}{5} = 33$$

$$S^2_{(\text{ave})} = \sqrt{\frac{(25-33)^2 + (27-33)^2 + \dots + (42-33)^2}{5}} = 7.823$$

$$P(\text{?o|No}) = \frac{1}{7.82 \times \sqrt{2\pi}} e^{-\frac{(20-33)^2}{2 \times (7.82)^2}} = 0.012$$

$$P(\text{low|No}) = \frac{1}{5} \quad P(\text{fair|No}) = \frac{2}{5}$$

$$P(\text{yes|No}) = \frac{1}{5}$$

$$P(\text{No|New}) = P(\text{?o|No}) * P(\text{low|No}) * P(\text{yes|No}) * P(\text{fair|No})$$

$$* P(\text{No})$$

$$= 0.012 * \frac{1}{5} * \frac{1}{5} * \frac{2}{5} * \frac{5}{9} = 6.85 * 10^{-5}$$

البحث \rightarrow New object \rightarrow يتنبئ \rightarrow class label
 \rightarrow Yes \Leftrightarrow Probabil \rightarrow الـ Prob

(4)

c) ? Values of Age all J_{age} Discrete

$$\text{متوسط} J_{age} = \frac{45 - 17}{2} = 14 \quad A \rightarrow B$$

[17, 31], [31, 45]

A: 17, 21, 25, 27, 28, 29

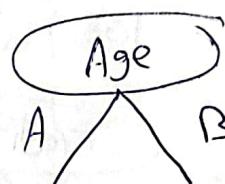
B: 31, 31, 31, 32, 36, 42, 43, 45

17, 21, 25, 27, 28, 29, 31,
31, 31, 32, 36, 42, 43,
45

لور طبقه هسته Decision Tree $\xrightarrow{\text{جذع}} 36$
Root لا ، سعر 4 جد Best Attr

| | |
|-----|---|
| No | 5 |
| Yes | 9 |

Data set



| | |
|-----|---|
| No | 3 |
| Yes | 3 |

| | |
|-----|---|
| No | 2 |
| Yes | 6 |

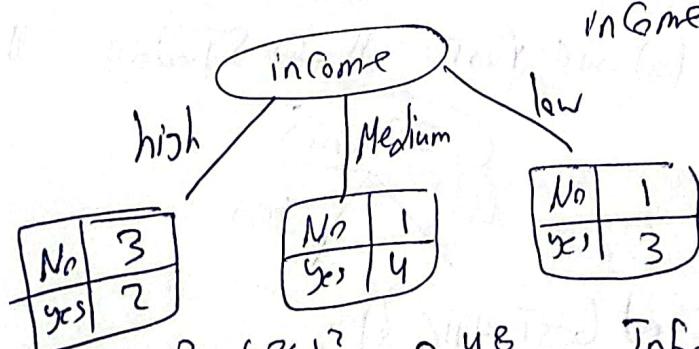
$$\text{Gini (Parent)} = 1 - \left(\frac{5}{14}\right)^2 - \left(\frac{9}{14}\right)^2 = \frac{45}{98}$$

$$\text{Gini (LeftT)} = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = \frac{1}{2}$$

$$\text{Gini (RightT)} = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = \frac{3}{8}$$

$$\text{Weighted Avg} = \frac{6}{14} \times \frac{1}{2} + \frac{8}{14} \times \frac{3}{8} = \frac{3}{7}$$

$$\text{Infor} = \frac{45}{98} - \frac{3}{7} = 0.0306 \quad \text{--- ①}$$



In Gini Job

$$\text{Gini (LeftT)} = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48$$

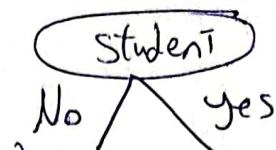
$$\text{Infor} = \frac{45}{98} - \frac{11}{28}$$

$$\sim (\text{MiddleP}) = 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{4}{5}\right)^2 = 0.32$$

$$= 0.6632 \quad \text{--- ②}$$

$$\checkmark (\text{low}) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0.375$$

$$\text{Weighted Avg} = \frac{5}{14} \times 0.48 + \frac{5}{14} \times 0.32 + \frac{4}{14} \times 0.375 = 0.392 = \frac{11}{28}$$



| | |
|-----|---|
| No | 4 |
| yes | 2 |

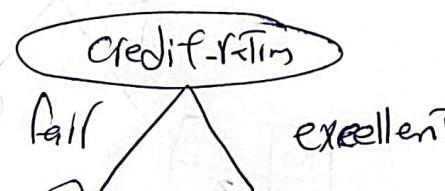
| | |
|-----|---|
| No | 1 |
| yes | 7 |

$$\text{Gini (left)} = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = \frac{4}{9}$$

$$\text{Gini (right)} = 1 - \left(\frac{1}{8}\right)^2 - \left(\frac{7}{8}\right)^2 = \frac{1}{32}$$

$$\text{Weighted AVG} = \frac{6}{14} * \frac{4}{9} + \frac{8}{14} * \frac{1}{32} = 0.3154 = \frac{53}{168}$$

$$\text{Info} = \frac{45}{98} - \frac{53}{168} = [0.1437] - ③$$



| | |
|-----|---|
| No | 3 |
| yes | 6 |

| | |
|-----|---|
| No | 2 |
| yes | 3 |

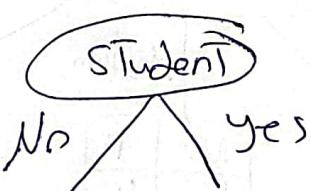
$$\text{Gini (left)} = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{6}{5}\right)^2 = \frac{4}{9}$$

$$\text{right} = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = \frac{12}{25}$$

$$\text{AVG} = \frac{9}{14} * \frac{4}{9} + \frac{5}{14} * \frac{12}{25} = \frac{16}{35}$$

$$\text{Info} = \frac{45}{98} - \frac{16}{35} = [2.04 * 10^{-3}] - ④$$

لـ ٤) اخيroot لـ Student لـ ٦) (جـ)



أ) best ATTR فـ وـ نـ وـ دـ اـ لـ ٦) (جـ)

٦)

| | |
|-----|---|
| No | 4 |
| Yes | 2 |

No



| | |
|-----|---|
| No | 3 |
| Yes | 0 |

| | |
|-----|---|
| No | 1 |
| Yes | 2 |

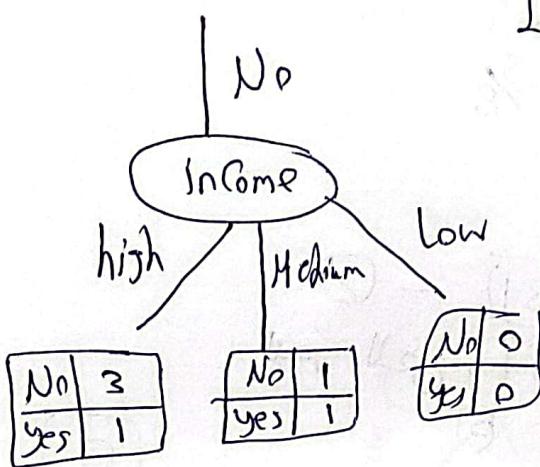
$$\text{Gini}(\text{left}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0$$

$$\text{Gini}(\text{right}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = \frac{4}{9}$$

$$\text{AVG} = \frac{3}{6} \times 0 + \frac{3}{6} \times \frac{4}{9} = \frac{2}{9}$$

$$\text{Gini}(\text{parent}) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2$$

$$I_m = \frac{4}{9} - \frac{2}{9} = \frac{2}{9} = 0.222 \quad \textcircled{1}$$



| | |
|-----|---|
| No | 3 |
| Yes | 1 |

| | |
|-----|---|
| No | 1 |
| Yes | 1 |

| | |
|-----|---|
| No | 0 |
| Yes | 0 |

$$\text{Gini}(\text{left}) = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{1}{6}\right)^2 = \frac{3}{8}$$

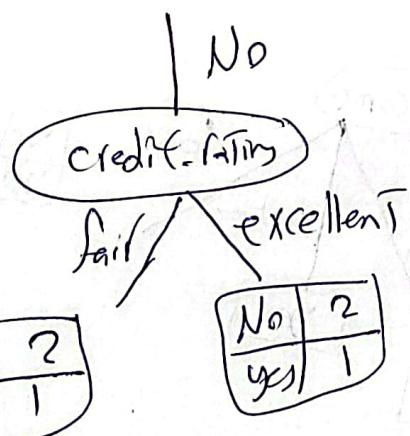
$$\text{Gini}(\text{middle}) = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{1}{6}\right)^2 = \frac{1}{3}$$

$$\text{right} = 0$$

$$\text{AVG} = \frac{4}{6} \times \frac{3}{8} + \frac{2}{6} \times \frac{1}{3} + 0 = \frac{5}{12}$$

$$I_m L = \frac{4}{9} - \frac{5}{12} = 0.027 \quad \textcircled{2}$$

7



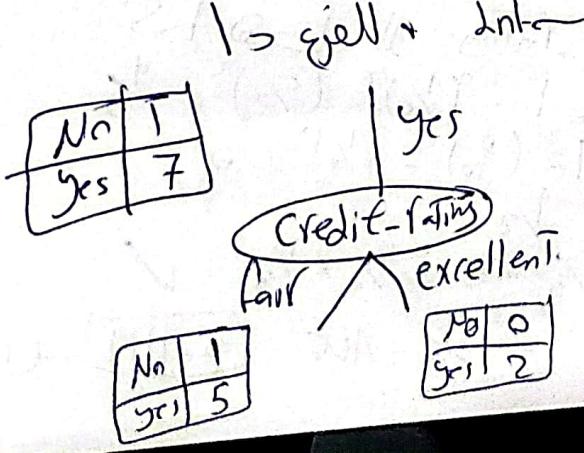
| | |
|-----|---|
| No | 2 |
| Yes | 1 |

| | |
|-----|---|
| No | 2 |
| Yes | 1 |

$$\text{Gini}(\text{left \& right}) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{1}{4}\right)^2$$

$$\text{AVG} = \left(\frac{3}{6} \times \frac{4}{9}\right) \times 2 = \frac{4}{9}$$

$$I_m F = \frac{4}{9} - \frac{4}{9} = 0 \quad \textcircled{3}$$



| | |
|-----|---|
| No | 1 |
| Yes | 5 |

| | |
|-----|---|
| No | 0 |
| Yes | 2 |

لبنان لا يدخل في المجموع

$$\text{Left} = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2 = \frac{5}{18}$$

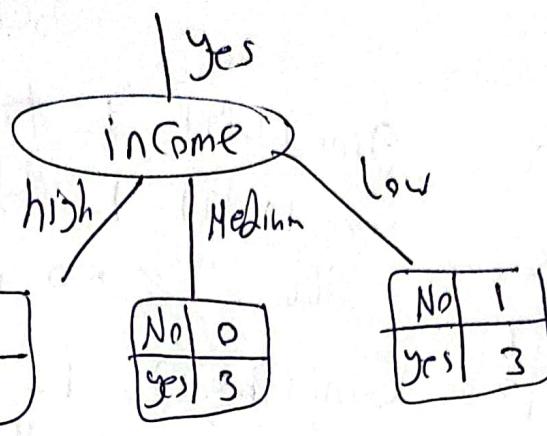
$$\text{Right} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{1}{2}$$

$$\text{AVG} = \frac{6}{12} + \frac{1}{12} + 0 = \frac{7}{12}$$

$$\text{Gini}(\text{parent}) = 1 - \left(\frac{1}{12}\right)^2 - \left(\frac{11}{12}\right)^2 = \frac{5}{24}$$

$$I_m L = \frac{1}{32} - \frac{5}{24} = \frac{7}{32} = 0.021875 \quad \textcircled{4}$$

| | |
|-----|---|
| No | 1 |
| Yes | 7 |



$$\text{Gini}(\text{leftT}) = 1 - \left(\frac{1}{8}\right)^2 - \left(\frac{7}{8}\right)^2 = 0$$

$$\text{Gini}(\text{Middle}) = 1 - \left(\frac{3}{8}\right)^2 - \left(\frac{5}{8}\right)^2 = 0$$

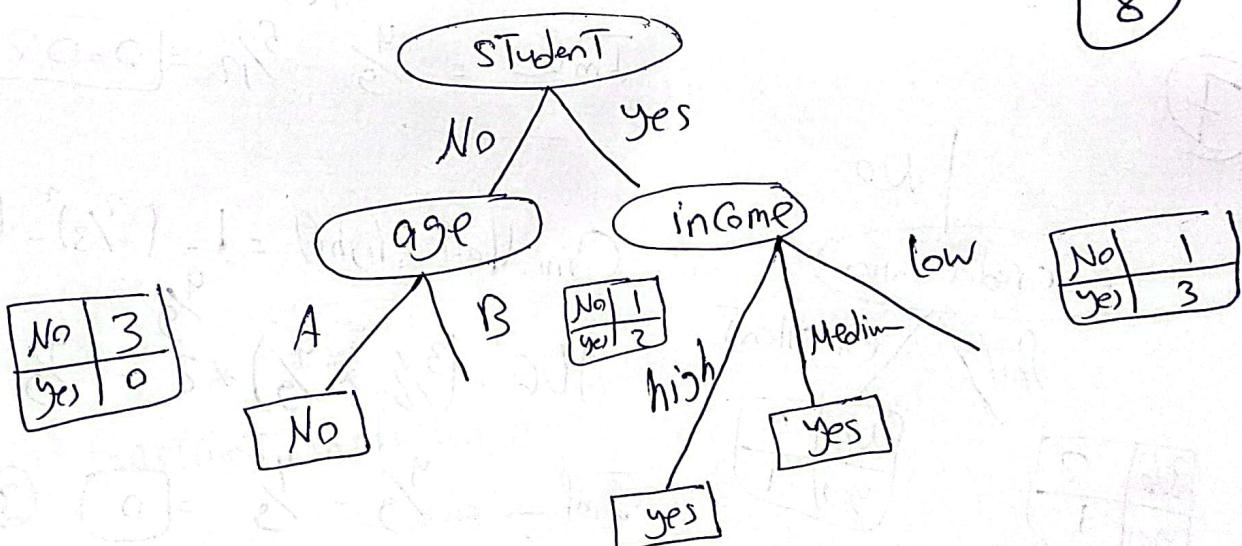
$$\text{Gini}(\text{RightT}) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = \frac{3}{8}$$

$$\text{AVG} = \frac{3}{8} * \frac{4}{8} = \frac{3}{16}$$

$$\text{InL} = \frac{1}{32} - \frac{3}{16} = \boxed{0.031} \rightarrow ②$$

so Income is a better feature (less InL)

8



so Age is a better feature for Credit Rating

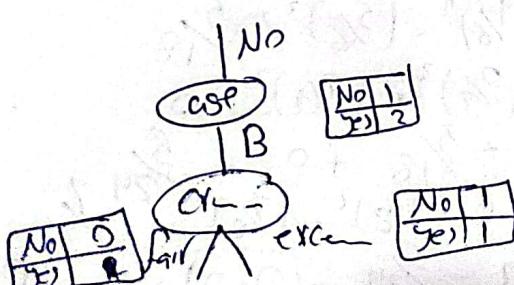
$$\text{Gini}(\text{Parent}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = \frac{4}{9}$$

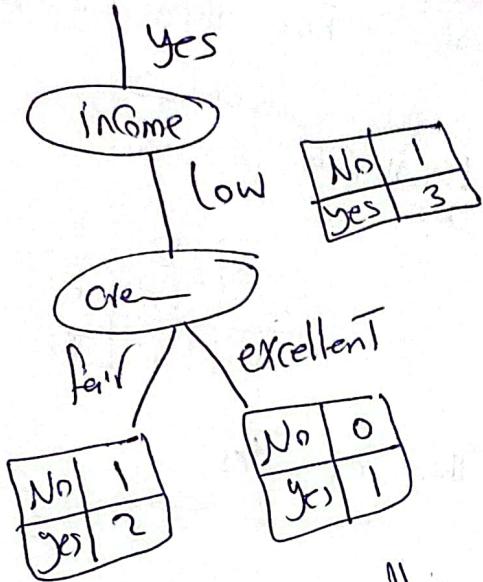
$$\text{Gini}(\text{LeftT}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0$$

$$\text{Gini}(\text{RightT}) = \frac{1}{2}$$

$$\text{AVG} = \frac{1}{3} * 0 + \frac{2}{3} * \frac{1}{2} = \checkmark$$

$$\text{InL} = \frac{4}{9} - \text{AVG} = \boxed{0.111} \rightarrow ①$$





$$Gini(\text{Parent}) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = \frac{3}{8}$$

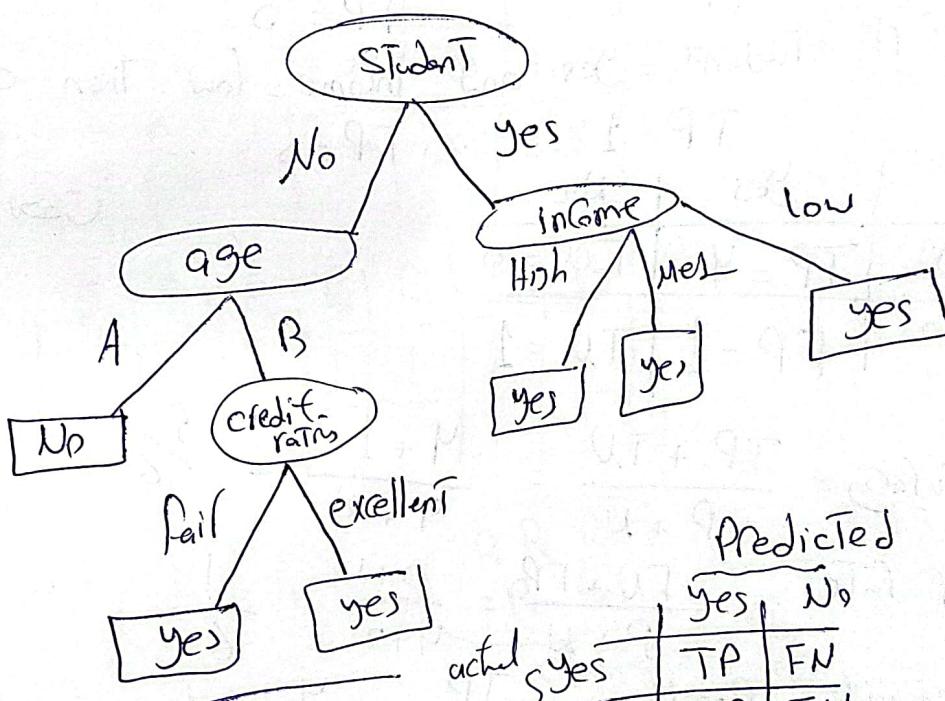
$$Gini(\text{Left}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = \frac{4}{9}$$

$$Gini(\text{Right}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0$$

$$AVG = \frac{3}{4} \times \frac{4}{9} = \frac{1}{3}$$

$$InL = \frac{3}{8} - \frac{1}{3} = 0.0411$$

لـ \Rightarrow الفـ لـ \Rightarrow Income \Rightarrow الـ \Rightarrow
IS هـ \Rightarrow Decision Tree



d) ii Construct Confusion Matrix \Rightarrow Tree Rule Rule \Rightarrow اختيار ونـ \Rightarrow

Rule1: if Student = No and age = A Then CL = No
 اـ لـ دـ حـ نـ دـ اـ بـ تـ \Rightarrow Predict
 if IS = No \Rightarrow CL \Rightarrow Predict
 yes or No \Rightarrow rule \Rightarrow rule \Rightarrow Test set \Rightarrow CL
 خـ اـ كـ دـ بـ نـ وـ الـ تـ بـ

$FN = 0$ $TN = 1$
 Predicted No
 Tree \Rightarrow No
 actual No
 Test set

(9)

* Pg 7

Rule2: if Student = No and Age = B and credit = Fair Then CL = yes

$$TP = 0 \quad FP = 1$$

As we see here
CL = If Rule 2 is not
yes

Rule3: if Student = No and Age = B and credit = excellent Then CL = yes

$$TP = 1 \quad FP = 0$$

| | | |
|-----|-----|----|
| | Yes | No |
| Yes | TP | FN |
| No | FP | TN |

Rule4: if Student = Yes and income = high Then CL = yes

$$TP = 1 \quad FP = 0$$

Rule5: if Student = Yes and income = Medium Then CL = yes

$$TP = 1 \quad FP = 0$$

Rule6: if Student = Yes and income = low Then CL = yes

$$TP = 1 \quad FP = 0$$

| | Yes | No |
|-----|--------|---------------|
| Yes | TP = 4 | <u>FN = 0</u> |
| No | FP = 1 | TN = 1 |

Now to calc.

$$\text{accuracy} = \frac{TP + TN}{P + N} = \frac{4 + 1}{4 + 2} = \frac{5}{6}$$

$$\text{error rate} = \frac{FN + FP}{P + N} = \frac{0 + 1}{4 + 2} = \frac{1}{6}$$

$$\text{Sensitivity} = \text{Recall} = \frac{TP}{P} = \frac{4}{4} = 1$$

$$\text{Specificity} = \frac{TN}{N} = \frac{1}{2}$$

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{4}{4 + 1} = \frac{4}{5}$$

(b)

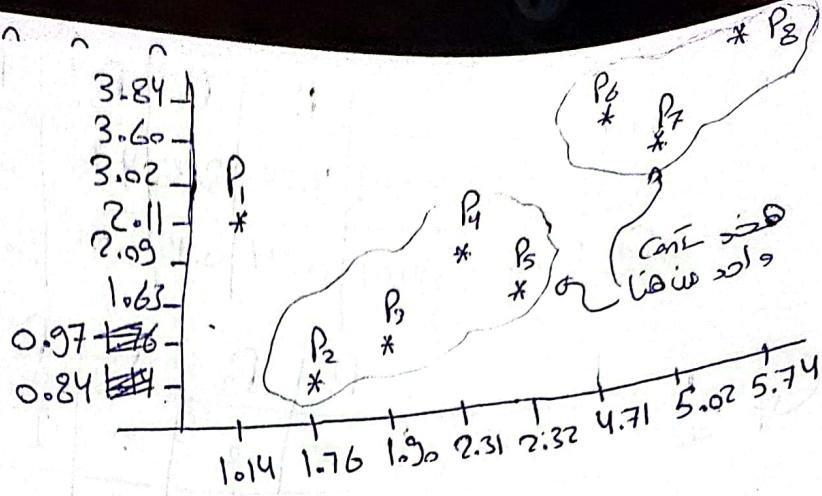
e) From rule 6

New object (?, low, yes, fair, ?)
is belong to class label yes

f) Similar

Q3

a) SCATTER PLOT



b) APPLY K-means Algo for K=2

$$C_1 = (5.02, 3.02), \quad C_2 = (2.31, 2.09)$$

D₀ ≈

| | P ₁ | P ₂ | P ₃ | P ₄ | P ₅ | P ₆ | P ₇ | P ₈ |
|-------------------------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| C ₁ = (5.02, 3.02) | 3.98 | 3.03 | 3.73 | 0 | 1.09 | 0.65 | 2.86 | 3.92 |
| C ₂ = (2.31, 2.09) | 1.17 | 0.46 | 1.09 | 7.86 | 3.85 | 2.83 | 0 | 1.36 |

G₀

| | P ₁ | P ₂ | P ₃ | P ₄ | P ₅ | P ₆ | P ₇ | P ₈ |
|-----------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| cluster 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| cluster 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

P₄, P₅, P₆

P₁, P₂, P₃
P₇, P₈

$$C_1 = \left(\frac{5.02 + 5.79 + 4.71}{3}, 3.48 \right)$$

$$\frac{3.02 + 3.84 + 3.60}{3}$$

iteration = ?

$$= \left(\frac{5.15}{5}, \frac{3.48}{5} \right)$$

$$C_2 = \left(\frac{1.04 + 2.32 + 1.90 + 2.31 + 1.76}{5}, \frac{2.11 + 1.63 + 0.97 + 2.09 + 0.84}{5} \right)$$

94

III

D₁)

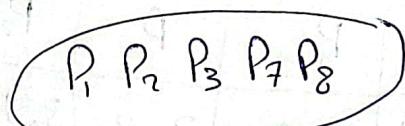
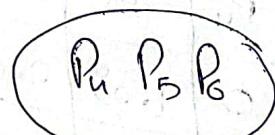
| | P ₁ | P ₂ | P ₃ | P ₄ | P ₅ | P ₆ | P ₇ | P ₈ |
|--------------------------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| C ₁ = (5.15, 3.48) | 4.23 | 3.38 | 4.10 | 0.47 | 0.69 | 3.46 | 3.16 | 4.29 |
| C ₂ = (1.88, 1.528) | 0.94 | 0.45 | 0.55 | 3.47 | 4.31 | 3.31 | 0.71 | 0.69 |

G₁ = 0

| | P ₁ | P ₂ | P ₃ | P ₄ | P ₅ | P ₆ | P ₇ | P ₈ |
|----------------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| cluster ₁ | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| cluster ₂ | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 |

نفس الـ G يتغير

STOP



i) Single link = أحياناً المعاقة بين لقرب نقطتين

cluster₂ و P₇ ، cluster₁ و P₆ و

$$Dis(P_6, P_7) = \sqrt{(4.71 - 2.31)^2 + (3.60 - 2.09)^2} = 2.83$$

ii) Complete link

المعاقة (بعد تقطيب)

P₆ و P₁ و

$$DB(P_1, P_6) = \sqrt{(1.14 - 4.71)^2 + (7.11 - 3.60)^2} = 3.86$$

iii) Centroid link

الكل هناCentroid أحياناً

$$C_1 = (5.15, 3.48)$$

$$C_2 = (1.88, 1.528)$$

$$DB(C_1, C_2) = \sqrt{(5.15 - 1.88)^2 + (3.48 - 1.528)^2} = 3.81$$

12

d) we choose P_1, P_2, P_3, P_4

| | P_1 | P_2 | P_3 | P_4 |
|-------|-------|-------|-------|-------|
| P_1 | 0 | | | |
| P_2 | 1.36 | 0 | | |
| P_3 | 1.37 | 0.78 | 0 | |
| P_4 | 3.98 | 3.03 | 3.73 | 0 |

0.78 \Rightarrow a new jpl
Merge P_3, P_2 into $J_{2,3}$

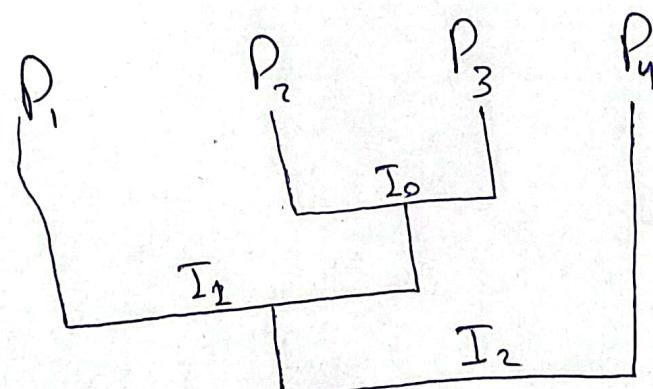
| | P_1 | (P_2, P_3) | P_4 |
|--------------|-------|--------------|-------|
| P_1 | 0 | | |
| (P_2, P_3) | 1.36 | 0 | |
| P_4 | 3.98 | 3.03 | 0 |

1.36 \Rightarrow a new jpl
Merge P_1 into $J_{1,2,3}$
 $P_1 \rightarrow (P_2, P_3)$

| | $(P_1, (P_2, P_3))$ | P_4 |
|---------------------|---------------------|-------|
| $(P_1, (P_2, P_3))$ | 0 | 0 |
| P_4 | 3.03 | 0 |

1st ↗

e)



$L=1$
 0.78
 $L=2$
 1.36
 $L=3$
 3.03

↳ cut $J_{2,3} \Rightarrow$

↳ result $\{I_2\}$

↳

P_1, P_2, P_3, P_4
 cluster

✗

13