

①

Linear regression

Q-1 From the following data find the best fit line using linear regression of sales on purchase.

Sale	91	97	108	121	67	124	51	73	111	57
Purchase	71	75	69	97	70	91	39	61	80	47

→ Given that

Sale	Purchase
91	71
97	75
108	69
121	97
67	70
124	91
51	39
73	61
111	80
57	47

We know that Best fit line using linear regression is:

$$y = ax + b$$

$$\text{where } a = \frac{n \sum xy - \sum x \sum y}{n(\sum x^2) - (\sum x)^2}$$

$$\text{and } b = \frac{\sum y - a \sum x}{n}$$

(2)

n	Sales (x)	Purchase (y)	xy	x <sup>2</sup>
1	91	71	6461	8281
2	97	75	7275	9409
3	108	69	7452	11664
4	121	97	11737	14641
5	67	70	4690	4489
6	124	91	11284	15376
7	51	39	1989	2601
8	73	61	4453	5329
9	111	80	8880	12321
10	57	47	2679	3249
	$\Sigma x = 900$	$\Sigma y = 700$	$\Sigma xy = 66900$	$\Sigma x^2 = 87360$

$n = 10$   
 $\Sigma x = 900$   
 $\Sigma y = 700$   
 $\Sigma xy = 66900$   
 $\Sigma x^2 = 87360$

For value of a

$$a = \frac{n \Sigma xy - \Sigma x \Sigma y}{n (\Sigma x^2) - (\Sigma x)^2}$$

$$= \frac{10 * 66900 - 900 * 700}{10 (87360) - (900)^2}$$

$$= \frac{1669000 - 630000}{873600 - 810000}$$

3

$$\frac{39000}{63600}$$

$$0.613$$

$$\therefore a = 0.613 \quad \text{--- (i)}$$

For value of b

$$b = \frac{\sum y - a \sum x}{n}$$

$$= \frac{(900 - 0.613(900))}{10}$$

$$= \frac{(900 - 551.07)}{10}$$

$$\therefore b = \frac{124.83}{10} = 12.483 \quad \text{--- (ii)}$$

$$y = ax + b \quad \text{--- (iii)}$$

From putting value of a and b in eqn (iii) from eqn (i) and (ii) we get

$$y = 0.613x + 12.483$$

is the best fit line using regression for given data of sales on purchase.

Q-2 Find the best line using linear regression of Math on English given the following data set also estimate the error.

English	20	60	55	45	75	35	25	90	10	50
Math	20	45	65	40	55	35	15	80	25	50

We know that

Best fit line using linear regression  $= y = ax + b$

where  $a = \frac{n \sum xy - \sum x \sum y}{n(\sum x^2) - (\sum x)^2}$

and  $b = \frac{\sum y - a \sum x}{n}$

Formula for error estimation  $\sigma_{est} = \sqrt{\frac{\sum y - y^2}{n}}$

n	English (x)	Math (y)	xy	x <sup>2</sup>
1	20	20	400	400
2	60	45	2700	3600
3	55	65	3575	3025
4	45	40	1800	2025
5	75	55	4125	5625
6	35	35	1225	1225
7	25	15	375	625
8	90	80	7200	8100
9	10	25	250	100
10	50	50	2500	2500

$n = 10$   
 $\sum x = 465$   
 $\sum y = 430$   
 $\sum xy = 24150$   
 $\sum x^2 = 27225$

(5)

Value of  $a$

$$a = \frac{n \sum xy - \sum x \sum y}{n(\sum x^2) - (\sum x)^2}$$

$$= \frac{10 * 24150 - 465 * 430}{10 * 27225 - (465)^2}$$

$$= \frac{241500 - 199950}{272250 - 216225}$$

$$= \frac{41550}{56025}$$

$$= \frac{41550}{56025}$$

$$= 0.742$$

$$\therefore a = 0.742 \quad \dots (1)$$

Value of  $b$

$$b = \frac{\sum y - a \sum x}{n}$$

$$= \frac{430 - 0.742 * 465}{10}$$

$$= \frac{(430 - 345.03)}{10}$$

$$\therefore b = 8.497 \quad \dots (2)$$

We know that  $y = ax + b$ , from (1) & (2)

$y$  will be

$$y = 0.742x + 8.497$$



(3)

For the error estimation, we know that

$$\sigma_{est} = \sqrt{\frac{\sum (y - y')^2}{n}}$$

n	Math (y)	y'	y - y'	(y - y') <sup>2</sup>
1	20	23.337	-3.337	11.135
2	45	53.017	-8.017	64.272
3	55	49.307	5.693	32.410
4	40	41.887	-1.887	3.550
5	55	64.147	-9.147	83.667
6	35	34.467	0.533	0.284
7	15	27.047	-12.047	145.130
8	80	75.277	4.723	22.306
9	25	15.917	9.083	82.500
10	30	45.597	4.403	19.386

$$\sum (y - y')^2 = 464.75$$

$$\therefore \sigma_{est} = \sqrt{\frac{464.75}{10}}$$

$$= \sqrt{46.475}$$

$$\therefore \sigma_{est} = 6.817$$

For best line using linear regression for math on english is  $y = 0.742x + 8.494$  and estimated error is 6.817.

(7)

Q-3

Find the least square regression line,  $y = ax + b$  for the following data. Also estimate the value of  $y$  when  $x = 10$ .

$x$	$y$
0	2
1	3
2	5
3	4
4	6

Formulae:

$$y = ax + b$$

$$\text{where } a = \frac{n \sum xy + \sum x + \sum y}{n(\sum x^2) - (\sum x)^2}$$

$$\text{and } b = \frac{\sum y - a \sum x}{n}$$

$x$	$y$	$xy$	$x^2$
0	2	0	0
1	3	3	1
2	5	10	4
3	4	12	9
4	6	24	16

$$\sum x = 10$$

$$\sum y = 20$$

$$\sum xy = 49$$

$$\sum x^2 = 30$$

(8)

Value of a

$$a = \frac{n \sum xy - \sum x \sum y}{n (\sum x^2) - (\sum x)^2}$$

$$= \frac{5(49) - 10 \times 20}{5 \times 30 - (10)^2}$$

$$= \frac{245 - 200}{150 - 100}$$

$$= \frac{45}{50}$$

$$\therefore a = 0.9$$

$$b = \frac{\sum y - a \sum x}{n} = \frac{1}{5} [20 - 0.9 \times 10] = 2.2$$

We know that

$$y = ax + b$$

$$y = 0.9x + 2.2$$

Given that  $x = 10$ ,  $y = ?$

$$y = 0.9(10) + 2.2$$

$$y = 11.2$$

$\therefore$  For given data value  $x = 10$  then  $y = 11.2$



(9)

# Bayesian Classifier

Q-1

Consider the following data of buying computer and classify a tuple  $X = (\text{age} \leq 30, \text{Income} = \text{medium}, \text{Student} = \text{yes}, \text{Credit rating} = \text{fair})$

Age	Income	Student	Credit rating	Buy computer
$\leq 30$	High	no	Fair	no
$\leq 30$	High	no	excellent	no
31-40	High	no	Fair	yes
$> 40$	Medium	no	Fair	yes
$> 40$	low	yes	Fair	yes
$> 40$	low	yes	excellent	no
31-40	low	yes	excellent	yes
$\leq 30$	medium	no	Fair	no
$\leq 30$	low	yes	Fair	yes
$> 40$	medium	yes	Fair	yes
$\leq 30$	medium	yes	excellent	yes
31-40	High	yes	Fair	yes
31-40	medium	no	excellent	yes
$> 40$	medium	no	excellent	no

$P(\text{C})$  :

$$P(\text{buy-computer} = \text{"yes"}) = 9/14 = 0.643$$

$$P(\text{buy-computer} = \text{"No"}) = 5/14 = 0.357$$

Formula :

Naive Bayes classifier :

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Given

$x = (\text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{Yes}, \text{Credit rating} = \text{Fair}, \text{buy-computer} = ?)$

Compute  $P(x/c_i)$  for each class

$$P(\text{age} \leq 30 \mid \text{buy-computer} = \text{Yes}) = 2/9 = 0.222$$

$$P(\text{age} \leq 30 \mid \text{buy-computer} = \text{No}) = 2/5 = 0.4$$

$$P(\text{income} = \text{medium} \mid \text{buy-computer} = \text{Yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} \mid \text{buy-computer} = \text{No}) = 2/5 = 0.4$$

$$P(\text{Student} = \text{Yes} \mid \text{buy-computer} = \text{Yes}) = 6/9 = 0.667$$

$$P(\text{Student} = \text{Yes} \mid \text{buy-computer} = \text{No}) = 1/5 = 0.2$$

$$P(\text{Credit rating} = \text{Fair} \mid \text{buy-computer} = \text{Yes}) = 6/9 = 0.667$$

$$P(\text{Credit rating} = \text{Fair} \mid \text{buy-computer} = \text{No}) = 2/5 = 0.4$$

Now

$P(x/c_i)$ :

$$P(x \mid \text{buy-computer} = \text{Yes})$$

$$= 0.222 \times 0.444 \times 0.667 \times 0.667$$

$$= 0.044$$

$$P(x \mid \text{buy-computer} = \text{No})$$

$$= 0.4 \times 0.4 \times 0.2 \times 0.4$$

$$= 0.019$$

$$P(x/c_i)P(c_i) \text{ For } P(\text{buy-computer} = \text{Yes})$$

$$P(x \mid \text{buy-computer} = \text{Yes}) \times P(\text{buy-computer} = \text{Yes})$$

$$= 0.044 \times 0.643 = 0.028$$

For  $\text{buy-computer} = \text{No}$

$$P(x \mid \text{buy-computer} = \text{No}) \times P(\text{buy-computer} = \text{No})$$

$$= 0.019 \times 0.357 = 0.007$$

$\therefore x$  belongs to Yes class  $\Rightarrow$  buy-computer.

Q-2 Consider the following data of approved loan  
 Class and classify a tuple  
 (Age=young, has\_job=false, own\_house=false, credit\_rating  
 =good, class=?)

ID	Age	has_job	own_house	Credit rating	Class
1	Young	F	F	Fair	No
2	Young	F	F	Good	No
3	Young	T	F	Good	Yes
4	Young	T	T	Fair	Yes
5	Young	F	F	Fair	No
6	Mid	F	F	Fair	No
7	Mid	F	F	Good	No
8	Mid	T	T	Good	Yes
9	Mid	F	T	Excellent	Yes
10	Mid	F	T	Excellent	Yes
11	Old	F	T	Excellent	Yes
12	Old	F	T	Good	Yes
13	Old	T	F	Good	Yes
14	Old	T	F	Excellent	Yes
15	Old	F	F	Fair	No

$P(\text{class})$  :

$P(\text{loan\_approved} = \text{Yes}) = 8/15 = 0.4$

$P(\text{loan\_approved} = \text{No}) = 7/15 = 0.6$

Formula :

Naive Bayes Classifier :

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Given that  
 $x$  (age = young, has job = false, own house = false, credit rating = good; class = ?)

Compute  $P(x | c_i)$  for each class

$$P(\text{age} = \text{young} | \text{approved} - \text{loan} = \text{yes}) = 2/6 = 0.33$$

$$P(\text{age} = \text{young} | \text{approved} - \text{loan} = \text{No}) = 4/9 = 0.44$$

$$P(\text{has job} = \text{false} | \text{approved} - \text{loan} = \text{yes}) = 4/6 = 0.67$$

$$P(\text{has job} = \text{false} | \text{approved} - \text{loan} = \text{No}) = 6/9 = 0.67$$

$$P(\text{own house} = \text{false} | \text{approved} - \text{loan} = \text{yes}) = 3/6 = 0.5$$

$$P(\text{own house} = \text{false} | \text{approved} - \text{loan} = \text{No}) = 6/9 = 0.67$$

$$P(\text{credit rating} = \text{good} | \text{approved} - \text{loan} = \text{yes}) = 4/6 = 0.67$$

$$P(\text{credit rating} = \text{good} | \text{approved} - \text{loan} = \text{No}) = 2/9 = 0.22$$

$$P(x | \text{loan approved} = \text{yes}) \\ = 0.33 \times 0.67 \times 0.5 \times 0.67 \\ = 0.074$$

$$P(x | \text{loan approved} = \text{No}) = 0.44 \times 0.67 \times 0.67 \times 0.22 \\ = 0.043$$

Now  $P(x | c_i) P(c_i)$

- for approved - loan = yes

$$P(x | \text{approved loan} = \text{yes}) \times P(\text{loan approved} = \text{yes}) \\ = 0.074 \times 0.4 = 0.029$$

$$P(x | \text{loan approved} = \text{No}) \times P(\text{loan approved} = \text{No}) \\ = 0.043 \times 0.6 \\ = 0.026$$

$\therefore x$  belong to class = yes as  $P(x | c_i)$  of loan approved - loan is higher than loan approved = No.



Q-3

With the help of the training sample given below using Bayes classification and can we believe that the patient having the given symptoms has Flu?

$x$  {chill=yes, rainy\_nose = No, Headache=Mild, Fever=yes, Flu=?}

Chills	Rainy-nose	Headache	Fever	Flu
Yes	No	Mild	Yes	No
Yes	Yes	no	no	Yes
Yes	No	Strong	Yes	Yes
No	Yes	mild	Yes	Yes
No	No	no	no	No
No	Yes	Strong	Yes	Yes
No	Yes	strong	no	No
Yes	Yes	mild	Yes	Yes

Formula :

Naive Bayes classifier

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$P(C_i)$

$$P(\text{Flu} = \text{Yes}) = 5/8 = 0.625$$

$$P(\text{Flu} = \text{No}) = 3/8 = 0.375$$

Calculate  $P(x|c_i)$  for each class  
 Given that

$x$  {chill=yes, rainy\_nose=no, headache=mild, fever=yes, Flu=?}

(14)

$$P(\text{chill} = \text{yes} \mid \text{flu} = \text{yes}) = 3/5 = 0.6$$

$$P(\text{chill} = \text{yes} \mid \text{flu} = \text{no}) = 1/3 = 0.33$$

$$P(\text{runny-nose} = \text{yes} \mid \text{flu} = \text{yes}) = 2/5 = 0.4$$

$$P(\text{runny-nose} = \text{no} \mid \text{flu} = \text{no}) = 2/3 = 0.66$$

$$P(\text{headache} = \text{mild} \mid \text{flu} = \text{yes}) = 2/5 = 0.4$$

$$P(\text{headache} = \text{mild} \mid \text{flu} = \text{no}) = 1/3 = 0.33$$

$$P(\text{fever} = \text{yes} \mid \text{flu} = \text{yes}) = 4/5 = 0.8$$

$$P(\text{fever} = \text{yes} \mid \text{flu} = \text{no}) = 1/3 = 0.33$$

$P(x \mid C)$ :

$$P(x \mid \text{flu} = \text{yes}) = 0.6 \times 0.2 \times 0.4 \times 0.8 = 0.0384$$

$$P(x \mid \text{flu} = \text{no}) = 0.33 \times 0.66 \times 0.33 \times 0.33 = 0.0227$$

$P(x \mid C) P(C)$

$$P(x \mid \text{flu} = \text{yes}) P(\text{flu} = \text{yes}) = 0.0384 \times 0.625 = 0.0239$$

$$P(x \mid \text{flu} = \text{no}) P(\text{flu} = \text{no}) = 0.0227 \times 0.375 = 0.0085$$

∴ With the given  $x$  symptoms patient belong to class yes of flu.



Q-1 Construct Decision Tree based on ID3 for the following training data.

Day	Outlook	Temp	Humidity	Wind	Decision
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Formula :

$$Gini\ index = 1 - \sum (P_i)^2$$

For  $i = 0$  to number of classes

Step 1: Find the Gini index for each feature.

Outlook	No. of Instances	Yes	No	Gini Index	Gini Index of Outlook
Sunny	5	2	3	0.48	0.342
Overcast	4	4	0	0	
Rain	5	3	2	0.48	

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

$$\text{Gini (Outlook = Overcast)} = 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 0$$

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

$$\text{Crini (Outlook)} = \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48$$

$$\therefore \text{Gini (Outlook)} = 0.342$$

Temperature					Gini For Temperature
	No of Instance	Yes	No	Gini Index	
Hot	4	2	2	0.5	0.440
Mild	6	4	2	0.445	
Cool	4	3	1	0.375	

$$\text{Gini (Temperature = hot)} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 0.5$$

$$\text{Gini (Temperature = mild)} = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.445$$

$$\text{Gini (Temperature = cool)} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

$$\text{Crini (Temperature)} = \frac{4}{14} \times 0.5 + \frac{6}{14} \times 0.445 + \frac{4}{14} \times 0.375$$

$$\therefore \text{Gini (Temperature)} = 0.440$$

Humidity					Gini Index for Humidity
	No. of Instances	Yes	No	Gini Index	
High	7	4	3	0.489	0.367
Normal	7	6	1	0.245	

$$\text{Gini (Humidity = high)} = 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = 0.489$$

$$\text{Gini (Humidity = normal)} = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0.245$$

$$\text{Crini (Humidity)} = \frac{7}{14} \times 0.489 + \frac{7}{14} \times 0.245$$

$$\therefore \text{Gini (Humidity)} = 0.367$$

Wind	No. of Instance	Yes	No	Gini Index	Gini Index for Wind
Strong	6	3	3	0.375	0.428
Weak	8	6	2	0.5	

$$\text{Gini}(\text{Wind} = \text{strong}) = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 0.375$$

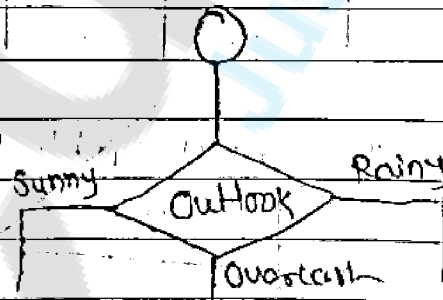
$$\text{Gini}(\text{Wind} = \text{weak}) = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.5$$

$$\text{Gini}(\text{Wind}) = \frac{6}{14} \times 0.375 + \frac{8}{14} \times 0.5$$

$$\therefore \text{Gini}(\text{Wind}) = 0.428$$

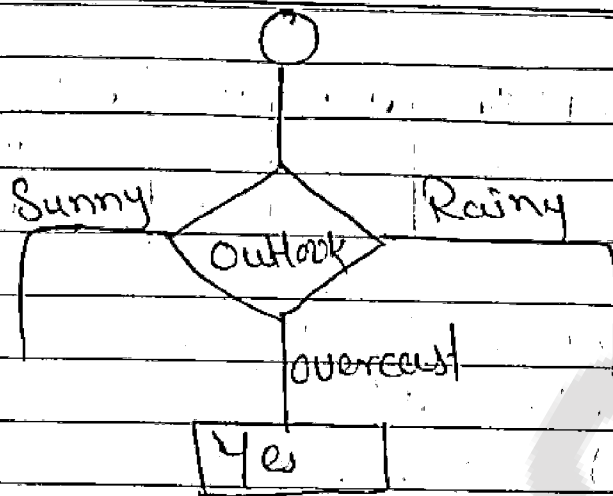
Step 2: Choose the lowest gini value feature as winner.

Feature	Gini Index
Outlook	0.342
Temperature	0.439
Humidity	0.367
Wind	0.428



From the given data, Overcast decision is always yes in any condition.

Hence the next feature will be overcast and decision will be "yes".



Step Now we have to find the decision Row, the Sunny outlook with respect to feature

Sunny Outlook

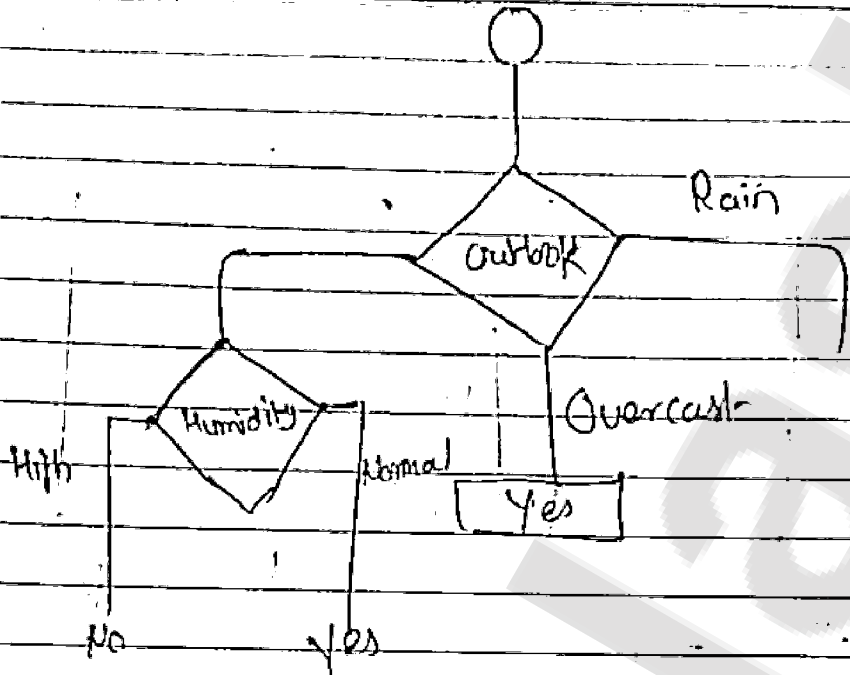
Temperature		No. of Instance	Yes	No	Gini of Instance	Gini of Sunny Instance
Cool	hot	2	0	2	0	
Hot	cool	1	1	0	0	0.2
	mild	2	1	1	0.5	
Humidity	High	3	0	3	0	0
	Normal	2	2	0	0	
Wind	Weak	3	1	2	2.66	0.466
	Strong	2	1	1	0.2	

Decision Row sunny outlook

Feature	Gini Index
Temperature	0.2
Humidity	0
Wind	0.466

Humidity win the sunny outlook

(14)



Rain outlook		yes	no	No. of Instances	Gini Index to Instance	Gini Index a Feature
Temperature	hot	0	0	0		
	cool	1	1	2	0.5	0.466
	mild	2	1	3	0.444	
Humidity	High	1	1	2	0.5	0.466
	Normal	2	1	3	0.444	
Wind	Weak	3	0	3	0	0
	Strong	0	2	2	0	

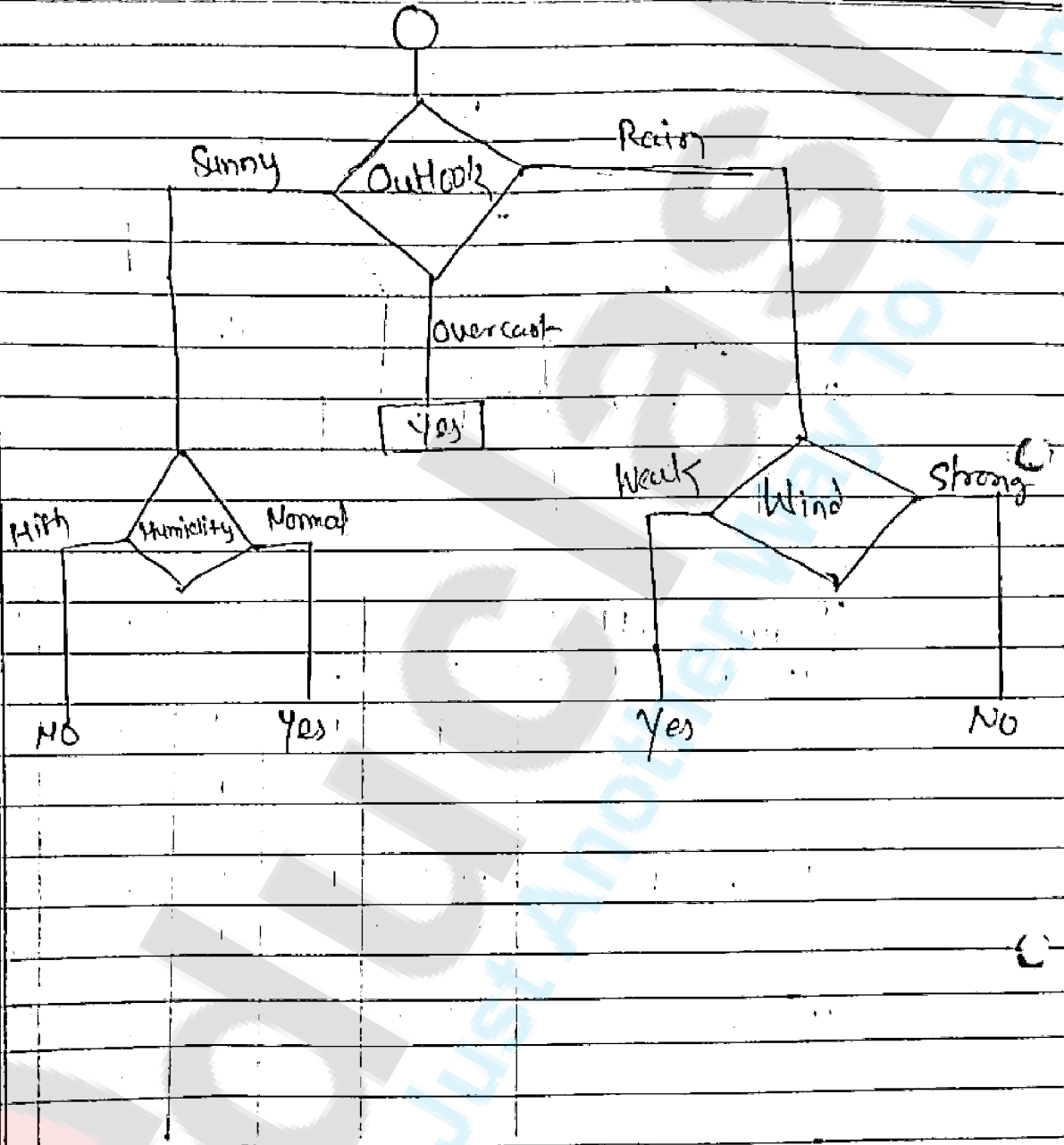
Decision For Rain outlook feature

Feature	Gini Index
Temperature	0.466
Humidity	0.466
Weak/Strong Wind	0

For rain outlooks the wind feature is important attribute.



2.0



2/5



Q-2

Generate the decision tree for following data set

No	Significant	Reorganized	Volume	Class
1	Yes	No	Small	Bad
2	Yes	Yes	large	Bad
3	No	Yes	medium	Bad
4	No	Yes	medium	Good
5	Yes	No	large	Good
6	No	No	large	Good

Formula :

$$\text{Gini Index} = 1 - \sum (P_i)^2$$

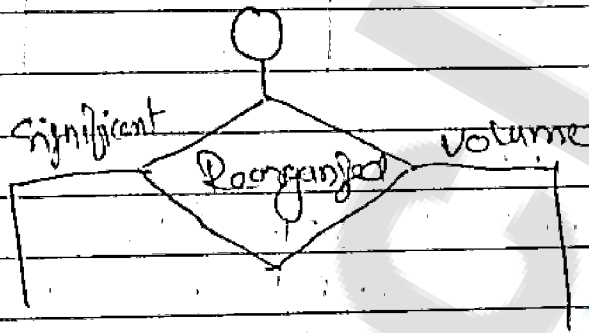
For p no number of classes

Find the gini index for each feature

Feature	No. of Instance	Good	Bad	Gini Index	Gini Feature
<b>Significant</b>					
Yes	3	1	2	0.444	0.444
No	3	2	1	0.444	
<b>Reorganized</b>					
Yes	2	0	2	0	0.25
No	4	3	1	0.375	
<b>Volume</b>					
Small	1	0	1	0	0.388
Mid	2	1	1	0.5	
Large	3	2	1	0.444	

Decision For Feature

Feature	Gini Index
Significant	0.444
Reorganized	0.25
Volume	0.388

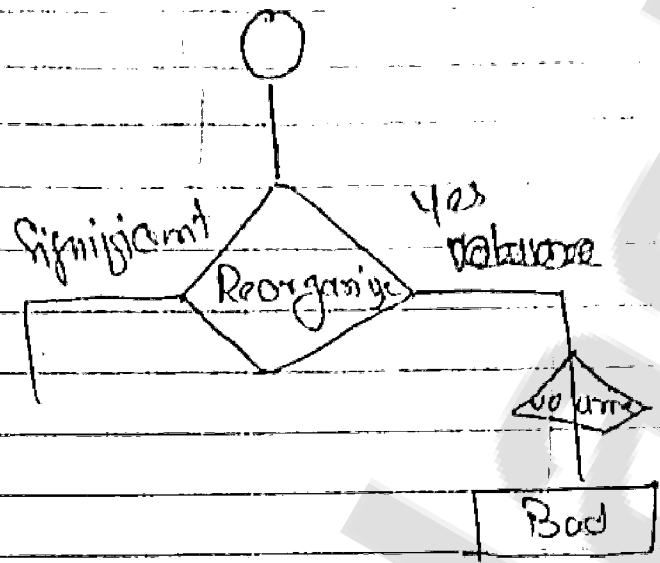


Yes reorganized

Feature	No. of	Instance	Good	Bad	Gini	Gini Index
Significant	0	3	1	2	0	0.444
	3					
	3					
Volume	1	3	0	1	0	0.388
	1		0	1	0	
	1		0	1	0	

It is clearly given that for yes reorganized division is bad.

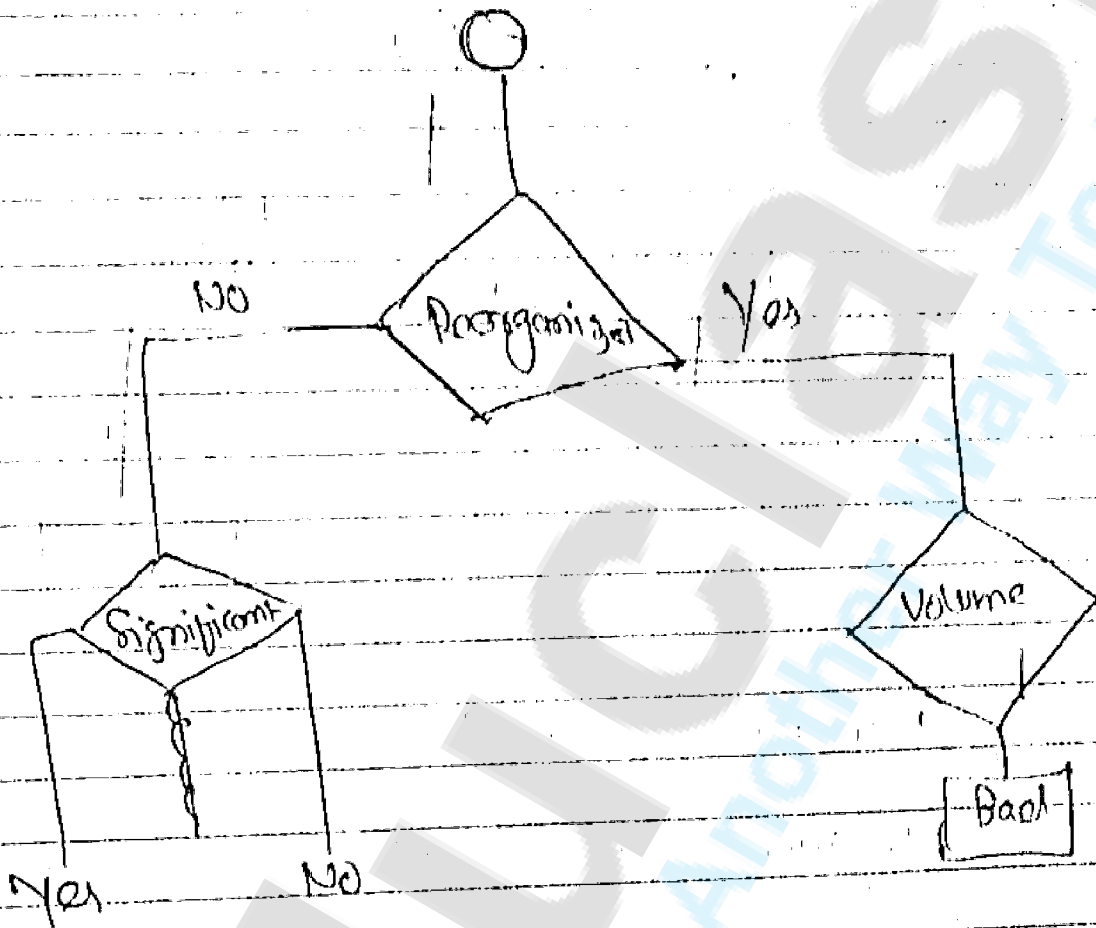
(22)



No reorganised

Feature	Nb. of Instance	Good	Bad	Final Instance	Final Feature
Significant					
Yes	2	1	1	0	0
No	2	2	0	0	
Volume					
Small	1	0	1	0	0
medium	1	0	1	0	
large	2	2	0	0	

Q. 81



05

Q-3 Construct a decision tree based on ID's for the following training data

	Chill	Rainy nose	Headache	Power	Flu
1	Yes	No	Mild	Yes	No
2	Yes	Yes	No	No	Yes
3	Yes	No	Strong	Yes	Yes
4	No	Yes	Mild	Yes	Yes
5	No	No	No	No	No
6	No	Yes	Strong	Yes	Yes
7	No	Yes	Strong	No	No
8	Yes	Yes	Mild	Yes	Yes

Formula

$$\text{Gini Index} = 1 - (\sum p_i)^2$$

Chill					Gini Index for Chill
	No. of Instance	Yes	No	Gini Index	
Yes	4	3	1	0.375	0.9375
No	4	2	2	0.5	

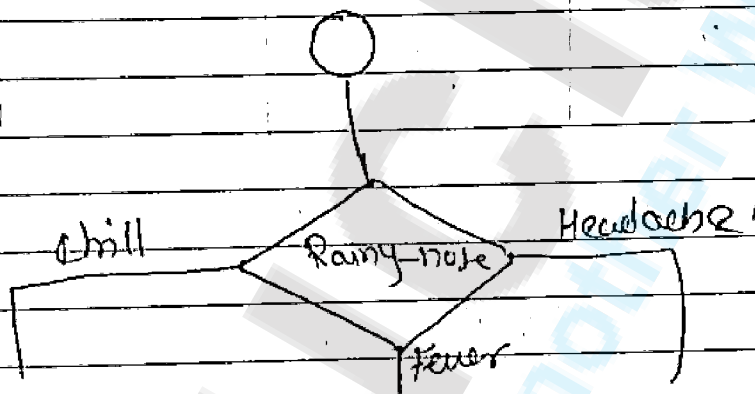
Rainy nose					
	No. of Instance	Yes	No	Gini Index	Gini Index
Yes	3	2	1	0.444	Rainy nose
No	5	2	3	0.32	0.366

Headache						
	No. of Instance	Yes	No	Gini Index	Gini for Headache	
No	2	1	1	0.5	0.458	
Mild	3	2	1	0.444		
Strong	3	2	1	0.444		



Fever		No. of Instance	Yes	No	Gini Index	Gini Index for Fever
Yes		5	4	1	0.32	
No		3	1	2	0.444	0.366

Feature	Gini Index
Chill	0.4735
Rainy-nose	0.366
Headache	0.458
Fever	0.366



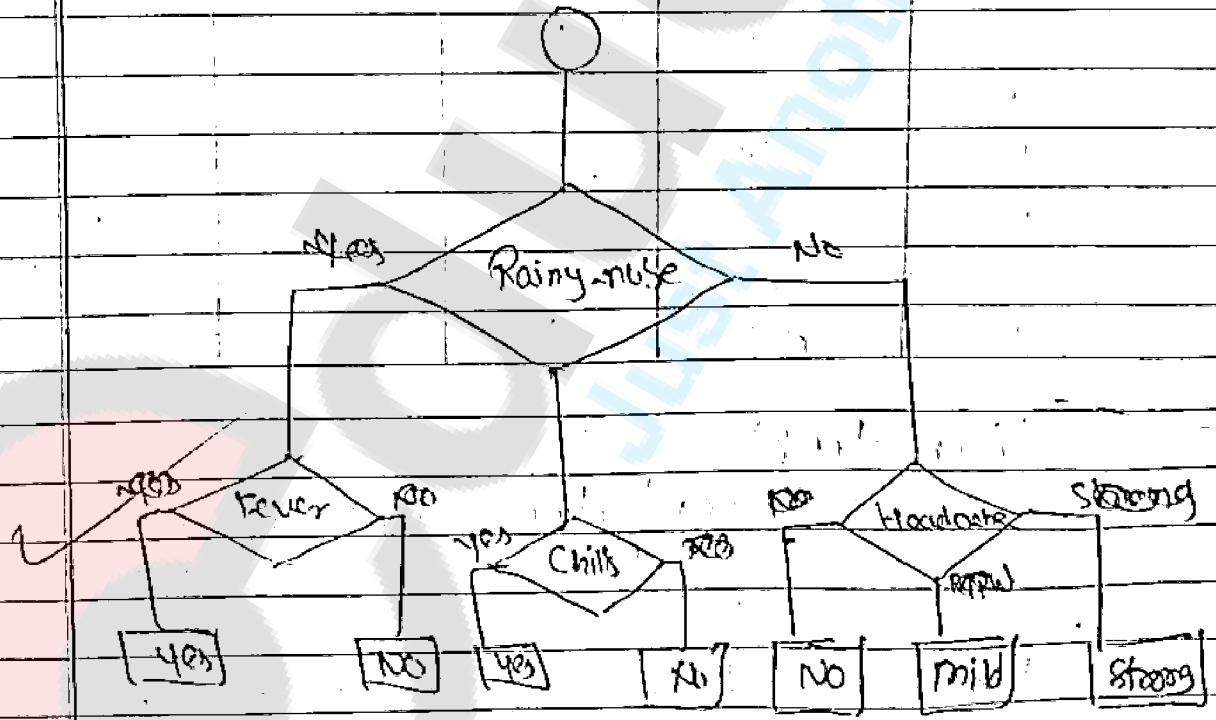
Rainy Nose = Yes

Feature	No. of Instance	Yes	No	Gini of balance	Gini of Feature
Chill					
Yes	2	2	0	0	0.266
No	3	2	1	0.444	
Headache					
No	1	1	0	0	0.2
Mild	2	2	0	0	
Strong	2	1	1	0.5	
Fever					
Yes	3	3	0	0	0.2
No	2	1	1	0.5	



Rainy nose = no

Feature	No. of Instance	Yes	No	Gini of Instance	Gini of Feature
Chill					
Yes	2	1	1	0.5	0.33
No	1	0	1	0	
Headache					
No	1	0	1	0	0
Mild	1	0	1	0	
Strong	1	1	0	0	
Fever					
Yes	2	1	1	0.5	0.33
No	1	0	1	0	



0.9%

Q-4 Construct a decision tree based on 103 for the following training set

ID	Age	Has job	Greenhouse	Excelling	Class
1	Young	F	F	Fair	No
2	Young	F	F	Good	No
3	Young	T	F	Good	Yes
4	Young	T	T	Fair	Yes
5	Young	F	F	Fair	No
6	Mid	F	F	Fair	No
7	Mid	F	F	Good	No
8	Mid	T	T	Good	Yes
9	Mid	F	T	Excellent	Yes
10	Mid	F	T	Excellent	Yes
11	Old	F	T	Excellent	Yes
12	old	F	T	Good	Yes
13	old	T	F	Good	Yes
14	old	T	F	Excellent	Yes
15	old	F	F	Fair	No

Formula:

$$Gini Index = 1 - \sum (P_i)^2$$

For 100 to number of classes

Find gini index for each feature.

Assignment No.: \_\_\_\_\_ Page No.: \_\_\_\_\_  
 Subject: \_\_\_\_\_  
 Roll No.: \_\_\_\_\_

28

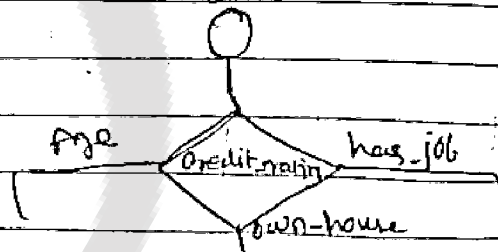
ICRD  
**Sterling**

INSTITUTE OF MANAGEMENT STUDIES, NERUL

Feature	No. of Instances	class		Gini Index of class	
		Yes	No		
Age	15				
Young	5	2	3	0.48	
Mid	5	3	2	0.48	0.426
Old	5	3	1	0.32	
has_job					
F	10	4	6	0.48	0.32
T	5	5	0	0	
Own-house					
F	9	3	6	0.44	0.264
T	6	6	0	0	
Credit-rating					
Fair	5	1	4	0.32	0.177
Good	6	4	2	0.44	
Excellent	4	4	0	0	

Feature	Gini Index
Age	0.426
has_job	0.32
Own-house	0.264
Credit-rating	0.177

As credit-rating has lowest Gini Index it is a main feature of given data.



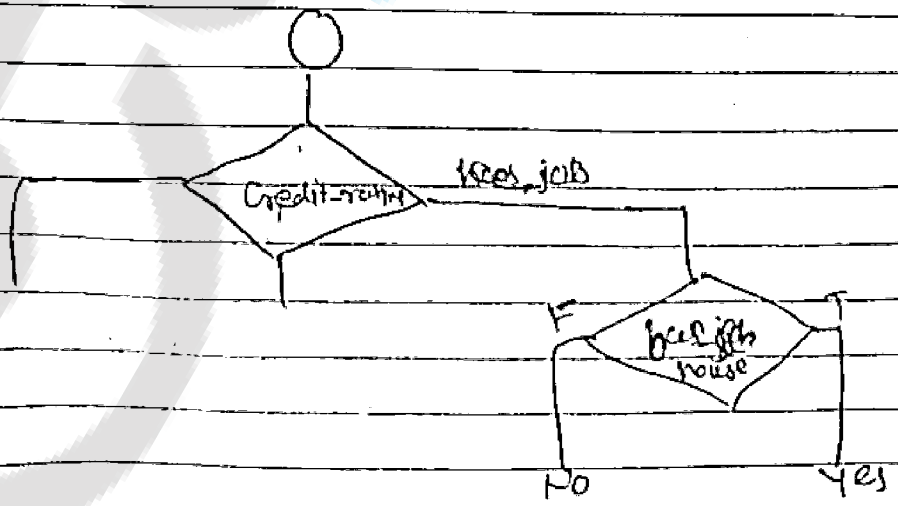
fair credit rating

Feature	No. of instances	Yes	No	Gini index for instance	Feature Gini Index
Age					
Young	3	1	2	0.44	
mid	1	0	1	0	
Old	1	0	1	0	0.269
has job					
False	4	0	4	0	
True	1	1	0	0	0
Own house					
F	4	0	4	0	0
T	1	1	0	0	

→ Decision

Feature	Gini Index
Age	0.269
has job	0
own-house	0

has job and own-house gain fair credit rating  
we will choose own house here

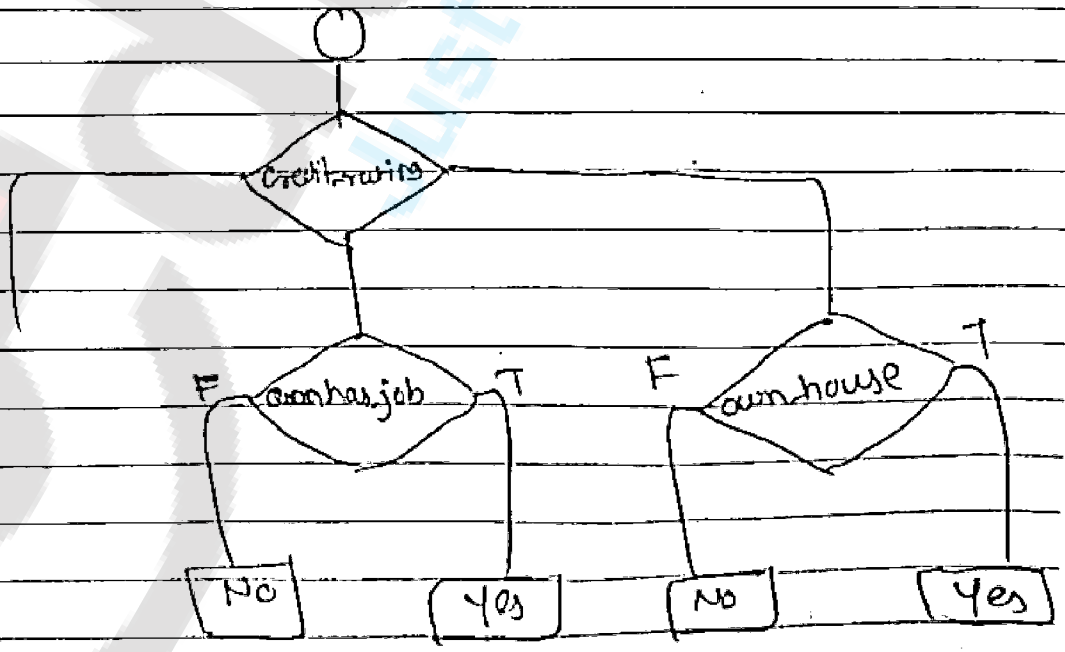


(31)

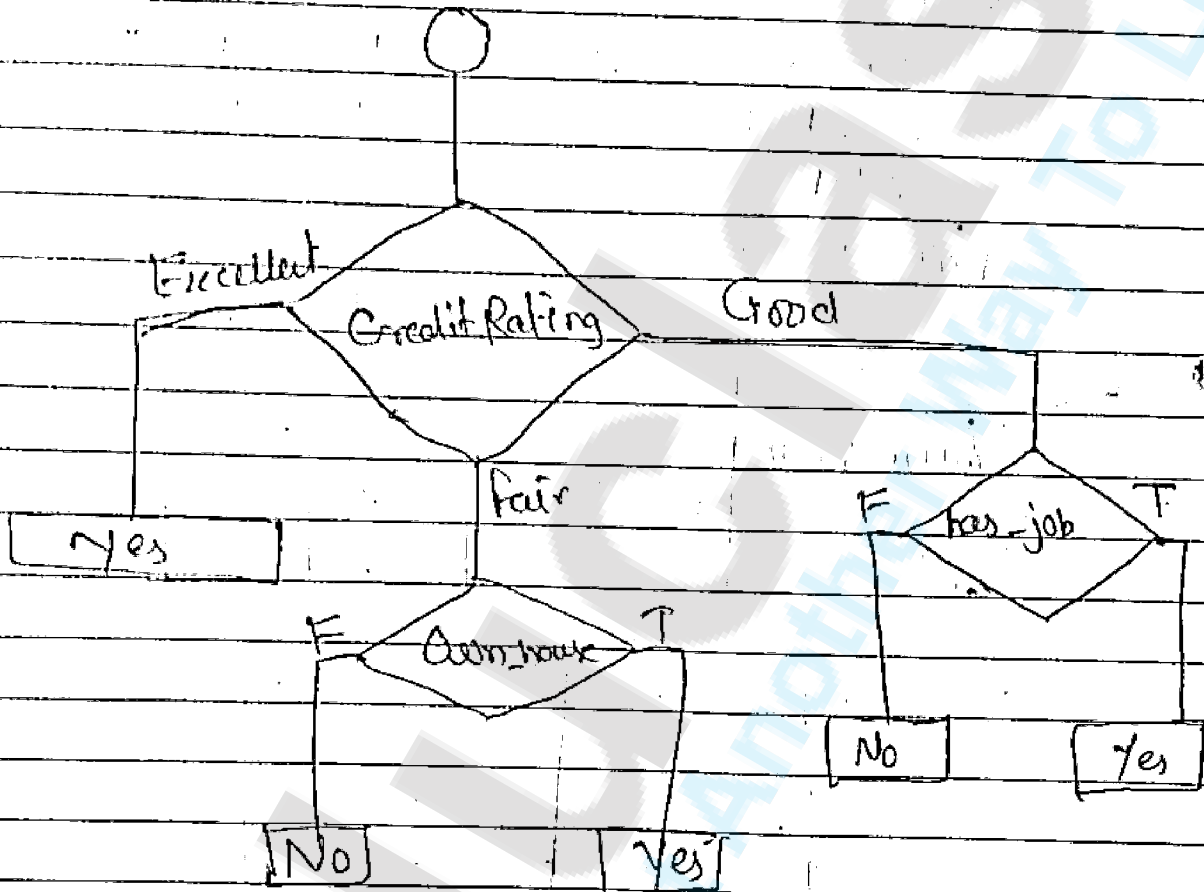
Good credit rating

Feature	No. of Instances	Yes	No	Gini Index	Gini Index for Feature
Age					
Young	2	1	1	0.5	
mid	2	1	1	0.5	0.33
old	2	2	0	0	
has job					
F	3	1	2	0.44	0.22
T	3	3	0	0	
own house					
F	4	2	2	0.5	0.33
T	2	2	0	0	

Feature	Gini Index
Age	0.33
has job	0.22
own house	0.33



For excellent credit rating it yes for each condition and/or feature.





P32

## K-means Clustering

Sterling

INSTITUTE OF MANAGEMENT STUDIES, NERUL

Q.1

What is clustering? Explain k-means clustering algorithm. Using k-means clustering, cluster the following data into two clusters and show each step.

{ 2, 4, 10, 12, 3, 20, 30, 11, 25 }

Clustering:- Clustering is a machine learning technique that involves the grouping of data points. Given a set of data points we can use a clustering algorithm to classify each data point into a specific group. In theory, data points are in the same group should have similar properties and/or features while data points in different groups should have less highly similar properties and/or features.

Clustering is a feature of unsupervised learning. It is a common technique for statistical data analysis used in many fields.

In data science, we can use clustering analysis to gain some valuable insight from our data by seeing what groups the data points fall into when we apply a clustering algorithm. There are 5 popular clustering algorithms:

- ① K-means clustering
- ② Mean-shift clustering
- ③ Density-Based Spatial Clustering of Applications with Noise (DBSCAN)
- ④ Expectation-Maximization (EM) clustering
- ⑤ Gaussian Mixture Model (GMM)
- ⑥ Agglomerative Hierarchical Clustering

34

k-mean clustering algorithm :-

- ① Randomly/given select cluster centres.
- ② Calculate the difference, between each data point and cluster centre
- ③ Assign data point to the cluster whose distance from the cluster centre is minimum from all the cluster centre.
- ④ Recalculate new cluster centre  

$$V_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (x_j)$$
- ⑤ Recalculate the distance between each data point and obtain new cluster centre.
- ⑥ If no data point was reassigned then stop otherwise repeat step 3.

k-mean clustering is pretty fast.

It has a linear complexity  $O(n)$ .

Q-8 Perform k-means clustering for given data set  
 $\{2, 3, 4, 10, 11, 12, 20, 25, 30\}$

→ Let  $m_1 = 4$  and  $m_2 = 12$   
 and  $K = 2$

$$K_1 = \{2, 3, 4\}$$

$$K_2 = \{10, 11, 12, 20, 25, 30\}$$

$$m_1 = \frac{2+3+4}{3} = 3$$

$$m_2 = \frac{10+11+12+20+25+30}{6} = 18$$

∴  $m_1 = 3$  and  $m_2 = 18$

→ New cluster

$$K_1 = \{2, 3, 4, 10\}$$

$$K_2 = \{11, 12, 20, 25, 30\}$$

$$\text{mean}(m_1) \text{ for } K_1 = 4.75 \approx 5$$

$$\text{mean}(m_2) \text{ for } K_2 = 19.6 \approx 20$$

→ New cluster

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

$$\text{mean}(m_1) \text{ for } K_1 = 7$$

$$\text{mean}(m_2) \text{ for } K_2 = 25$$

→ New cluster

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

$$\text{mean}(m_1) \text{ for } K_1 = 7$$

$$\text{mean}(m_2) \text{ for } K_2 = 25$$

Final cluster will be

Cluster ( $k_1$ )	Cluster ( $k_2$ )
2	20
3	25
4	30
10	
11	
12	

Step used in this numerical:

Step 1: Choose random cluster and initialize.

Step 2: Take a mean value.

Step 3: Find nearest number to mean and put in cluster.

Step 4: Repeat step (2) and step (3) until you get same mean value and/or you covered all the datapoint.

Step 5: Stop.

Q-21 Explain k-means clustering algorithm, perform k-means clustering using 1 Euclidean distance measure for the given data set ( $k=2$ ).

A	1.0	1.5	3.0	5.0	3.5	4.5	3.5
B	1.0	2.0	4.0	7.0	5.0	5.0	4.5

k-means clustering algorithm:

- ① Randomly select cluster centres.
- ② Calculate the difference between <sup>each</sup> data point and cluster centre.
- ③ Assign data point to the cluster whose distance from the cluster centre is minimum from all the cluster centre.

④ Recalculate the cluster centre.

$$V_i = \left( \frac{1}{U_i} \right) \sum_{j=1}^n (x_j^i)$$

⑤ Calculate the distance between each data point and obtain new cluster centre.

⑥ If no data point was reassigned then stop otherwise repeat step 6.



Obj ID	X	Y
1	1	1
2	1.5	2
3	3	4
4	5	7
5	8.5	5
6	4.5	5
7	8.5	4.5

$C_1 = \text{Cluster 1} = \{1, 2, 3\}$

Centroid for  $C_1 = [(1, 1), (1.25, 1.5), (2, 1.75)]$

$C_2 = \text{Cluster 2} = \{4, 5, 6, 7\}$

Centroid for  $C_2 = [(5, 7), (4.25, 5), (4.375, 5.5), (3.937, 5)]$

Formula:

$$\text{Euclidean Distance} = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$$

Initialize the cluster

Obj	X	Y	Centroid value
$C_1$	1	1	1
$C_2$	5	7	

→ For Row 2 using Euclidean formula  
for row we have to find Euclidean distance with respect to both cluster 1 and cluster 2

$$\begin{aligned} \text{ED for } C_1 &= \sqrt{(1.5 - 1)^2 + (2 - 1)^2} \\ &= \sqrt{1.25} \\ &= 1.11 \end{aligned}$$

$$\begin{aligned} \text{ED for } C_2 &= \sqrt{(1.5 - 5)^2 + (2 - 7)^2} \\ &= \sqrt{30.25} \\ &= 5.5 \end{aligned}$$

(29)

Row 2 goes in  $C_1$ .

$$\therefore \text{New centroid} = \left( \frac{1+1.5}{2}, \frac{1+2}{2} \right)$$

Value of  $C_1$

$$\text{New centroid value} = [1.25, 1.5]$$

→ For Row 3, using Euclidean Formula

$$\begin{aligned} \text{ED for } C_1 &= \sqrt{(3-1.25)^2 + (4-1.5)^2} \\ &= \sqrt{3.0625 + 6.25} = \sqrt{9.3125} \\ &= 3.05 \end{aligned}$$

$$\begin{aligned} \text{ED for } C_2 &= \sqrt{(3-5)^2 + (4-7)^2} \\ &= 3.60 \end{aligned}$$

Row 3 goes in cluster ( $C_1$ ).

$$\begin{aligned} \text{New centroid for } C_1 &= \left( \frac{1.5+3}{2}, \frac{2+4}{2} \right) \\ &= \left( \frac{4.5}{2}, \frac{6}{2} \right) \end{aligned}$$

$$\therefore \text{New Centroid for } C_1 = (2.25, 3)$$

→ For Row 5 using Euclidean formula

$$\begin{aligned} \text{ED for } C_1 &= \sqrt{(3.5-2.25)^2 + (5-3)^2} \\ &= 2.160 \end{aligned}$$

$$\begin{aligned} \text{ED for } C_2 &= \sqrt{(3.5-4)^2 + (5-5.5)^2} \\ &= 2.500 \end{aligned}$$

Row 5 goes in cluster 2.

$$\text{New centroid for } C_2 = \left( \frac{3.5+5}{2}, \frac{7+5}{2} \right)$$

$$\therefore \text{Centroid for } C_2 = (4.25, 6)$$

→ For Row 6, using Euclidean distance formula

$$\text{ED for } C_1 = \sqrt{(4.5 - 2.14)^2 + (5 - 2.75)^2} = 3.260$$

$$\text{ED for } C_2 = \sqrt{(4.5 - 4.25)^2 + (5 - 6)^2} = 1.030$$

Row 6 goes in cluster 2.

$$\text{New centroid} = \left( \frac{4.25+4.5}{2}, \frac{6+5}{2} \right)$$

$$\text{New centroid for } C_2 = (4.375, 5.5)$$

→ For Row 7, using euclidean formula

$$\text{ED for } C_1 = \sqrt{(3.5 - 2.14)^2 + (4.5 - 2.75)^2} = 2.216$$

$$\text{ED for } C_2 = \sqrt{(3.5 - 4.375)^2 + (4.5 - 5.5)^2} = 1.328$$

Row 7 goes in cluster 2.

$$\text{New centroid for } C_2 = \left( \frac{4.375+3.5}{2}, \frac{5.5+4.5}{2} \right) = (3.9375, 5)$$

Assignment No.: \_\_\_\_\_ Page No.: \_\_\_\_\_

Subject: \_\_\_\_\_

Roll No.: \_\_\_\_\_

41

NCRD's

**Sterling**

INSTITUTE OF MANAGEMENT STUDIES, NERUL

ID	X	Y	cluster 1		Centroid	ID	cluster 2		Centroid
			ID	X	( $\bar{x}, \bar{y}$ )		ID	X	( $\bar{x}, \bar{y}$ )
1	1	1	1	1	(1,1)	4	5	7	(5,7)
2	1.5	2	2	1.5	(1.25, 1.5)	5	3.5	5	(4.25, 5)
3	3	4	3	3	(2.14, 2.75)	6	4.5	5	(4.375, 5.5)
4	5	7				7	3.5	4.5	(3.937, 5)
5	3.5	5							
6	4.5	5							
7	3.5	4.5							

*John*  
v. good.

Just Another Nerul

Q-3 Explain the k-means clustering and write the algorithm. Find two clusters of the following data set.

	i	$x_1$	$x_2$
A	1	1	1
B		1	0
C		0	2
D		2	4
E		3	5

k-means clustering algorithm:

(1) Randomly select cluster centres.

(2) Calculate the difference between each data point and cluster centre.

(3) Assign data point to the cluster whose distance from the cluster centre is minimum from all the cluster centre.

(4) Recalculate the cluster centre.

$$V_i = (1/n_i) \sum_{j=1}^{n_i} (x_{ij})$$

(5) Recalculate the distance between each data point and obtain new cluster centre.

(6) If no data point was reclassified then stop otherwise repeat step 3.



43

Given that  $k = 2$

Let A and C are cluster 1 and cluster 2 respectively.

Given data set	X	Y
A	1	1
B	1	0
C	0	2
D	2	4
E	3	5

Formula :

$$\text{Euclidean Distance}^{(ED)} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Ans  $\rightarrow$

Cluster  $k_1 = \{A, B\}$

Cluster  $k_2 = \{C, D, E\}$

Centroid for  $k_1 = \{(1, 1), (1, 0.5)\}$

Centroid for  $k_2 = \{(0, 2), (1, 3)\}$

For Row B:

$$ED \text{ for cluster } k_1 = \sqrt{(1-1)^2 + (0-1)^2}$$

$$= 1$$

$$ED \text{ for } k_2 = \sqrt{(1-0)^2 + (2-2)^2}$$

$$= 1$$

Row B goes in cluster  $k_1$

New Centroid for  $k_1 = \left(\frac{1+1}{2}, \frac{1+0}{2}\right)$

Centroid for  $k_1 = (1, 0.5)$

