

Name : Bilal Malik
 SUBJECT : ML
 Roll No : 34-BSCS-17
 Submitted To : Prof. YESHA ATM

Question No. 1

Solutions:

$$P(\text{Yes}) = \frac{1}{2}$$

$$P(\text{No}) = \frac{1}{2}$$

Total # of Records

$$\begin{array}{l} P(\text{Yes}) \\ (5) \end{array}$$

$$\begin{array}{l} P(\text{No}) \\ (5) \end{array}$$

Origin

DOMESTIC (6)

IMPORTED (4)

Yes (3)

No (3)

Yes (2)

No (2)

$$P(\text{domestic/Yes}) = \frac{3}{5}$$

$$P(\text{domestic/No}) = \frac{3}{5}$$

$$P(\text{imported/Yes}) = \frac{2}{5}$$

$$P(\text{imported/No}) = \frac{2}{5}$$

Type

Sports (6)

SUV (4)

Yes (4)

No (2)

Yes (1)

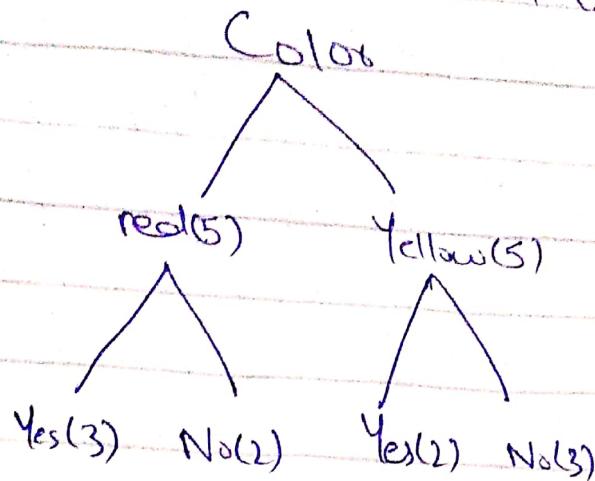
No (3)

$$P(\text{Sports}/\text{Yes}) = 4/5$$

$$P(\text{Sports}/\text{No}) = 2/5$$

$$P(\text{SUV}/\text{Yes}) = 1/5$$

$$P(\text{SUV}/\text{No}) = 3/5$$



$$P(\text{RED}/\text{Yes}) = 3/5$$

$$P(\text{RED}/\text{No}) = 2/5$$

$$P(\text{Yellow}/\text{Yes}) = 2/5$$

$$P(\text{Yellow}/\text{No}) = 3/5$$

LET

$$X = \{\text{red, SUV, Domestic}\}$$

$$P(X/\text{Yes}) = P(\text{Yes}) \times P(\text{RED}/\text{Yes}) \times P(\text{SUV}/\text{Yes}) \\ \times P(\text{Domestic}/\text{Yes})$$

$$= 1/2 \times 3/5 \times 1/5 \times 3/5$$

$$= 0.036$$

$$P(X/\text{No}) = P(\text{No}) \times P(\text{RED}/\text{No}) \times P(\text{SUV}/\text{No})$$

$$\times P(\text{Domestic}/\text{No})$$

$$= 1/2 \times 2/5 \times 3/5 \times 3/5$$

$$= 0.072$$

$$\text{As } 0.072 > 0.036$$

So THE PROBABILITY OF CAR

RANCE STOLEN IS NO.

QUESTION NO# 2.

DECISION TREE

$$P(\text{YES}) = 9 \quad P(\text{NO}) = 5$$

$$\text{Entropy} (S) = -\frac{9}{14} \log_2 \left[\frac{9}{14} \right] = \frac{5}{14} \log_2$$

$$\frac{5}{14} = .940$$

⁹⁺⁵

Calculating entropy for each attribute

Outlook	Yes	No	Entropy
SUNNY	2	3	.971
Rain	3	2	.971
OVERCAST	4	0	0

Using formula

$$E(\text{Outlook} = \text{SUNNY}) = -\frac{2}{5} \log \left(\frac{2}{5} \right) = \frac{3}{5} \log \left(\frac{3}{5} \right) = .971$$

Similarly calculate others

Information

CALCULATING Average Entropy

$$H(\text{Outlook}) = \frac{3+2}{9+5} \times .971 + \frac{2+3}{9+5} \times .971 + \frac{4}{14}$$

$$= .693$$

Now calculating gain

$G_{\text{ain}} = \text{Entropy}(S) = I$ (ATTRIBUTE)

$$\text{Entropy}(S) = .940$$

$$G_{\text{AIN}}(\text{Outlook}) = .940 - .693 = .247$$

For TEMPERATURE

Temp	Yes	No	Entropy
Hot	2	2	1
Mild	4	2	.918
Cool	3	1	.918

Avg INFORMATION ENTROPY

$$I(\text{Temp}) = \frac{2+2}{9+5} \times 1 + \frac{4+2}{9+5} \times .918 + \frac{3+1}{9+5} \times .81 \\ = .911$$

$$\text{Entropy}(S) = .940$$

$$G_{\text{AIN}} = .940 - .911 = .029$$

For HUMIDITY

Humidity	Yes	No	Entropy
High	3	4	.995
Normal	6	1	.591

Avg INFORMATION ENTROPY

$$I(\text{Humidity}) = \frac{3}{4} \times -0.985 + \frac{1}{4} \times 0.591$$

$\frac{3}{4}$

$\frac{1}{4}$

$$E(S) = -0.940$$

$$GAIN(\text{Humidity}) = 0.940 - 0.783 = 0.152$$

For Wind

Wind	Yes	No	Entropy
Strong	3	3	1
Weak	6	2	0.811

AVERAGE Info Entropy

$$I(\text{Windy}) = \frac{3}{4} \times 1 + \frac{1}{4} \times 0.811 = 0.8912$$

$\frac{3}{4}$

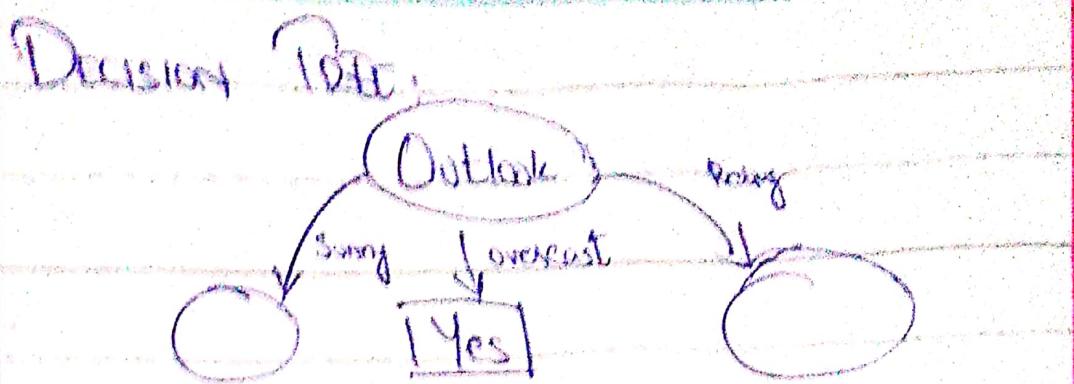
$\frac{1}{4}$

$$E(S) = -0.940$$

$$GAIN(\text{Wind}) = 0.940 - 0.892 = 0.048$$

SINCE outlook has biggest gain, it will be our root node

Attribute	Gain
Outlook	0.247
Temp	0.029
Humidity	0.152
Wind	0.048



Here further splitting is needed

Calculate entropy for sunny & rainy
of Outlook

For these 2 values,

$$P(\text{Yes}) = \frac{2}{5}, P(\text{No}) = \frac{3}{5}$$

$$\text{TOTAL} = 5$$

$$\text{Entropy}(S) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$E(\text{SUNNY}) = .971$$

For Humidity.

HUMIDITY	YES	NO	ENTROPY
HIGH	0	3	0
NORMAL	2	0	0
Avg Info Entropy			

$$I(\text{Humidity}) = 0$$

$$C_{\text{RAIN}} = .971$$

(4)

For Wind

Wind	Yrs	No	Entropy
Strong	1	1	1
Weak	1	2	.92

$$I(Windy) = .951 \quad GAIN = .020$$

For Temperature

Temp	Yes	No	Entropy
Cool	1	0	0
Hot	0	2	0
Mild	1	1	1

$$I(Temp) = .4$$

$$GAIN = .571$$

Now

ATTRIBUTES

GAINS

Temperature

.571

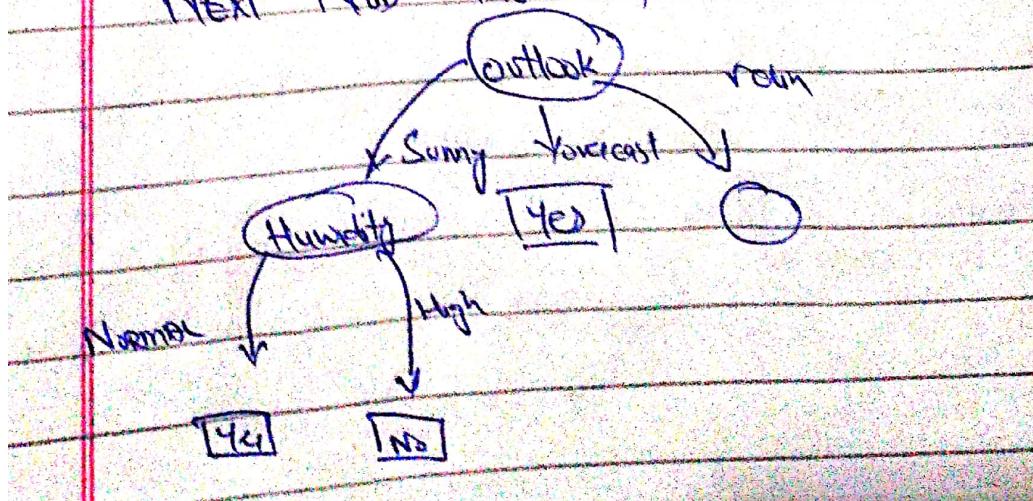
Humidity

.971

Windy

.02

Next No Altimetry



Now for rain, further splitting is necessary

$$P=3, N=2$$

Total = 5

$$\text{Entropy (S)}_{\text{rainy}} = - \frac{3}{3+2} \log_2 \frac{3}{3+2} - \frac{2}{3+2} \log_2 \frac{2}{3+2}$$
$$= .971$$

Calculating Entropy For Each ATTRIBUTE
For Humidity

Avg Info Entropy.

$$I(\text{Humidity}) = .951$$

$$GAIN = .020$$

For Wind

Avg Info Entropy

$$I(\text{Wind}) = 0$$

$$GAIN = .971$$

For Temp,

$$I(\text{Temp}) = .951$$

$$GAIN = .020$$

ATTRIBUTES

GAIN

HUMIDITY

$$.02$$

WIND

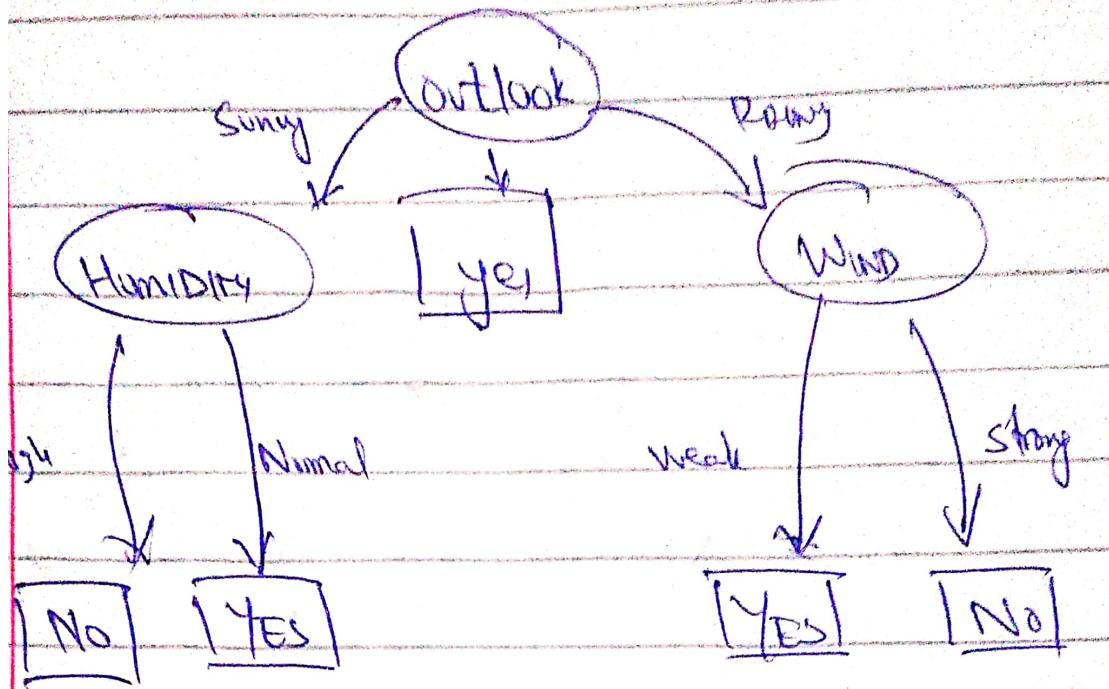
$$.971$$

TEMPERATURE

$$.02$$

Next Node WIND

FINAL DECISION TREE:



QUESTION: 3 (a)

K-MEAN ALGORITHM:

(2, 3, 4, 10, 11, 12, 20, 25 & 30) K=2

$$C_1 = 4 \quad \text{and} \quad C_2 = 12$$

By finding the distance of dataset

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

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

$$C_1 = 2+3+4/3 \approx 5$$

$$C_2 = 10+11+12+20+25+30/6 \approx 18$$

Now for Iteration #2

$$C_1 = 3 \quad \text{and} \quad C_2 = 18$$

Now As for dataset

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

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

$$C_1 = 2+3+4+10/4 = 4.75 \approx 5$$

$$C_2 = 11+12+20+25+30/5 = 19.6 \approx 20$$

For Iteration #3

$$C_1 = 5 \quad \text{and} \quad C_2 = 20$$

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

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

$$C_1 = 2+3+4+10+11+12/6 \approx 7$$

$$20+25+30/3 \approx 25$$

For iteration # 4

$$c_1 = 7$$

$$c_2 = 25$$

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

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

As K_1 , E_1 , K_2 of this iteration
are same as iteration # 4

$$c_1 = 7$$

$$c_2 = 25$$

So The cluster : 18

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

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

Q No 3 (B)

$(2,3), (5,6), (8,7), (1,4), (2,2), (6,7)$
 $(3,4), (8,6)$

Take any 2 points

① $c_1 = (2,3) \quad c_2 = (5,6)$

② Assign values in 2 clusters

Distance formula: $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$

x	y	distance $c_1(2,3)$	distance $c_2(5,6)$	Select cluster
2	3	0	4.24	c_1
5	6	4.24	0	c_2
8	7	7.21	3.16	c_2
1	4	1.41	4.47	c_1
2	2	1	5	c_1
6	7	5.65	1.41	c_2
3	4	1.41	2.82	c_1
8	6	6.70	3	c_2

③ Check which one in c_1 & c_2 :

$c_1 = (2,3), (1,4), (2,2), (3,4)$

Mean (x, y) = $(x_1 + x_2 + x_3 + x_4)/4, (y_1 + y_2 + y_3 + y_4)/4$

= $(2+1+2+3)/4, (3+4+2+4)/4$

$(2, 3, 3.5)$

$$C_2 = (5, 6)(8, 7)(6, 7)(8, 6)$$

$$\text{Mean}(x, y) = \frac{5+8+6+8}{4}, \quad \frac{6+7+7+6}{4}$$

$$C_2 = (6.75, 6.5)$$

Mean is different repeating steps

X	Y	distance	distance	selected cluster
2	3	= .25	5.9	C ₁
5	6	4.069	1.820	C ₂
8	7	7.075	1.346	C ₂
1	4	1.25	6.269	C ₁
2	2	1.25	6.54	C ₁
6	7	5.482	.90	C ₂
3	4	1.25	4.5	C ₁
8	6	6.005	1.346	C ₁

$$C_1 = (2, 3)(1, 4)(2, 2)(3, 4)$$

$$\text{Mean}(x, y) = \frac{2+1+2+3}{4}, \quad y = \frac{3+4+2+4}{4}$$

$$C_1 = (2, 3.25)$$

$$C_2 = (5, 6)(8, 7)(6, 7)(8, 6)$$

$$\text{Mean}(x, y) = \frac{5+8+6+8}{4}, \quad y = \frac{6+7+7+6}{4}$$

$$C_2 = (6.75, 6.5)$$

There is some, so stop it.

QUESTION No 4

	B/A	B/N	case	case	case	B/A	B/N
a	0.00	2.83	4.24	4.12	2.24	6.4	4.12
b	2.83	0	5.16	3.61	3	6.71	6.08
c	4.24	5.10	0	2.24	2.24	2.24	2.24
d	4.12	3.61	2.24	0	2.0	3.16	4.24
e	2.24	3	2.24	2	0	4.24	3.16
f	6.40	6.71	2.24	3.16	4.24	0	4.0
g	4.12	6.08	2.24	4.24	3.16	4	0

So, here we have

Min-points = 3

ϵ -Neighborhood = 2.25

CORE POINTS: c, d, e

BORDER POINTS: a, f, g

NOISE POINTS: b