

Q.2 Consider following dataset with min. support count = 60%. Find all frequent itemset using apriori algo & also generate strong association rules if min. confidence = 0.50%.

Transaction ID	Item Bought
T1	{M, D, N, K, E, Y}
T2	{D, O, N, K, E, Y}
T3	{M, A, K, E}
T4	{M, U, G, K, Y}
T5	{C, O, O, K, I, E}

Hint: O bought 4 times in total, but it occurs in just 3 transactions.

$$\begin{aligned}
 \text{Min Support Count} &= \text{min. support threshold} \\
 &\quad + \text{total no. of trans} \\
 &= \frac{60}{100} \times 5 \\
 &= 3
 \end{aligned}$$

$$\text{Min Support} = 60\%.$$

Database .

TID ITEM .

T1 MONKEY

T2 DONKEY

T3 MAKE

T4 MUCKY

T5 COOKIE

C1

Scan 0

Itemset Sup.

M	3	C	A
---	---	---	---

O	3	I	1
---	---	---	---

N	2		
---	---	--	--

K	5		
---	---	--	--

E	4		
---	---	--	--

Y	3		
---	---	--	--

D	1		
---	---	--	--

A	1		
---	---	--	--

U	1		
---	---	--	--

L1	itemset	Sup	C2	itemset
M	3			{M, OY}
O	3			{M, KY}
K	5			{M, EY}
E	4			{M, Y}
Y	3			{O, KY}

C2

L2	itemset	sup	itemset	sup	Scan	itemset
{M, KY}	3	{M, OY}	1		D	{K, EY}
{O, KY}	3	{M, KY}	3			{K, Y}
{O, EY}	3	{M, EY}	2			{E, Y}
{K, EY}	4	{M, Y}	2			
{K, Y}	3	{O, KY}	3			
		{O, EY}	3			
		{O, Y}	2			
		{K, EY}	4			
		{K, Y}	3			
		{E, Y}	2			

L3	itemset	sup	Scan D	itemset	sup
{O, K, EY}				{O, K, EY}	3
{K, E, Y}				{K, E, Y}	2

Frequent items \Rightarrow O, K, E

$$L = \{O, K, E\}$$

$$S = \{OY, KY, EY, \{O, K\}, \{O, E\}, \{K, E\}\}$$

$S \rightarrow (L-S)$	Support	Confidence	Confidence (%)
$\{OY \rightarrow \{K, E\}$	3	3/3	100
$\{KY \rightarrow \{O, E\}$	3	3/3	100
$\{EY \rightarrow \{O, K\}$	3	3/4	75
$\{O, K\} \rightarrow \{E\}$	3	3/3	100
$\{O, E\} \rightarrow \{KY\}$	3	3/3	100
$\{K, E\} \rightarrow \{OY\}$	3	3/4	75

∴ all the ass. rules are more than 50%.

∴ Strong association rules are

$$\{OY \rightarrow \{K, E\}$$

$$\{O, K\} \rightarrow \{EY\}$$

$$\{O, E\} \rightarrow \{KY\}$$

$$\{KY \rightarrow \{O, E\}$$

$$\{EY \rightarrow \{O, K\}$$

$$\{K, E\} \rightarrow \{OY\}$$

CLASSIFICATION TECHNIQUES

Decision Tree Classification

Ex Create classification model using decision tree:

S.R.NO	INCOME	AGE	OWN HOUSE
1	Very high	Young	Yes
2	High	Medium	Yes
3	Low	Young	Rented
4	High	Medium	Yes
5	Very high	Medium	Yes
6	Medium	Young	Yes
7	High	Old	Yes
8	Medium	Medium	Rented
9	Low	Medium	Rented
10	Low	Old	Rented
11	High	Young	Yes
12	Medium	Old	Rented

$$\text{Soln}^{\circ}- \text{No. of records} = 12$$

$$\text{Yes} = 7, \text{Rented} = 5$$

$$\text{Entropy (Own house)} = -\sum_{i=1}^{12} p_i \log_2 p_i$$

$$-0.1372 = 0.158$$

$$E(F) = E(7/12, 5/12)$$

$$= E(0.58, 0.41)$$

$$= -(0.58 \log_2 0.58) - (0.41 \log_2 0.41)$$

$$= 0.98$$

$$0.45582$$

$$0.5273$$

decrease in entropy after data
 DATE

target class

↓ Income

cex

Step 1 for $E(\text{Income})$ we have $E(T, X) = \sum P(c) E(c)$.

INCOME	YES	RENTED	TOTAL
Very High	2	6	2
High	4	0	4
Low	0	3	3
Medium	1	2	<u>3</u> <u>12</u>

$$E(\text{Own house, Income}) = p(vh) * E(vh) + \\ p(h) * E(h) + \\ p(l) * E(l) + \\ p(m) * E(m),$$

$$= \left[\left(\frac{2}{12} \right) * E\left(\frac{2}{2}, \frac{0}{2}\right) \right] +$$

$$\left[\left(\frac{4}{12} \right) * E\left(\frac{4}{4}, \frac{0}{4}\right) \right] + \left[\left(\frac{3}{12} \right) * E\left(\frac{0}{3}, \frac{3}{3}\right) \right]$$

$$+ \left[\left(\frac{3}{12} \right) * E\left(\frac{1}{3}, \frac{2}{3}\right) \right].$$

$$= 0 + 0 + 0 + \left[\left(\frac{3}{12} \right) E\left(\frac{1}{3}, \frac{2}{3}\right) \right]$$

$$= 0.25 * \left[(0.33 \log_2 0.33) - (0.67 \log_2 0.67) \right]$$

$$= 0.25 * 0.92$$

$$= 0.23$$

Step 2: For $E(\text{Age})$ we have $E(T, x) = \sum_{c \in X} p(c) E(c)$

AGE	YES	RENTED	TOTAL
Young	3	1	4
Medium	3	2	5
Old	2	2	3
			12

$E(\text{Own house}, \text{Age})$

$$= p(y) * E(y) + p(m) * E(m) + p(o) * E(o)$$

$$= \left[\left(\frac{4}{12} \right) * E\left(\frac{3}{4}, \frac{1}{4}\right) \right] + \left[\left(\frac{5}{12} \right) * E\left(\frac{3}{5}, \frac{2}{5}\right) \right]$$

$$+ \left[\left(\frac{3}{12} \right) * E\left(\frac{1}{3}, \frac{2}{3}\right) \right]$$

$$= \left[0.333 * \left[-(0.75 \log_2 0.75) - (0.25 \log_2 0.25) \right] \right]$$

$$+ \left[0.416 * \left[-(0.6 \log_2 0.6) - (0.4 \log_2 0.4) \right] \right]$$

$$+ \left[0.25 * \left[-(0.33 \log_2 0.33) - (0.66 \log_2 0.66) \right] \right]$$

$$= 0.333 * 0.3$$

$$= \underline{\underline{0.9}}$$

$$\text{Gain}(O, I) = E(O) - E(O, I)$$

$$= 0.98 - 0.23$$

$$= 0.75$$

—

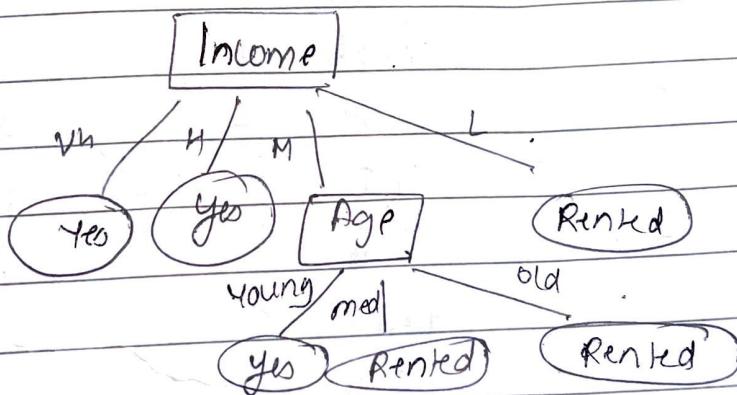
$$\text{Gain}(T, X) = E(T) - E(T, X)$$

$$\text{Gain}(O, A) = E(O) - E(O, A)$$

$$= 0.98 - 0.9$$

$$= \underline{\underline{0.08}}$$

Income attribute has highest gain, so used as a decision attribute in the root node.



Explain K-means clustering & solve the following with $K=3$ {2, 3, 6, 8, 9, 12, 15, 18, 22}

Solution :-

K-means clustering is an unsupervised learning algorithm.

K-means algorithm helps to group the unlabeled dataset into different clusters.

Here; K defines ~~as~~ the number of predefined clusters that need to be created in the process.

It is a centroid based algorithm, where each cluster is associated with a centroid.

The aim of the algorithm is to minimize the sum of distances between the data point and clustering their corresponding cluster.

The algorithm takes the unlabeled dataset as input, divides the dataset into K number of clusters, & repeats the process until it does not find the best clusters.

Given :-

Dataset : {2, 3, 6, 8, 9, 12, 15, 18, 22}

$K=3$

Taking the three mean,

$$M_1 = 2; M_2 = 3; M_3 = 6$$

Dividing the dataset as per the mean.

$$K_1 = \{2\}$$

$$K_2 = \{3\}$$

$$K_3 = \{6, 8, 9, 12, 15, 18, 22\}$$

Now calculating new mean as per the cluster created.

$$M_1 = 2; M_2 = 3; M_3 = \frac{(6+8+9+12+15+18+22)}{7} = 12.85$$

Dividing the dataset as per the newly created mean;

$$K_1 = \{2\}$$

$$K_2 = \{3, 6, 8, 9\}$$

$$K_3 = \{12, 15, 18, 22\}$$

Now, we will calculate a new mean as per new clusters:

$$M_1 = 2; M_2 = \frac{(3+6+8+9)}{4} = 6.5 \quad M_3 = \frac{12+15+18+22}{4}$$

Dividing the dataset as per the newly created mean;

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

$$K_2 = \{6, 8, 9\}$$

$$K_3 = \{12, 15, 18, 22\}$$

Now, we will calculate a new mean as new cluster:

$$M_1 = 2.5; M_2 = \frac{6+8+9}{3} = 7.67, M_3 = \frac{(12+15+18+22)}{4}$$

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

$$K_2 = \{6, 8, 9, 12\}$$

$$K_3 = \{15, 18, 22\}$$

Now, we will calculate our new mean again
and open the new clusters

$$M_1 = 2.5 \quad M_2 = 8.75 \quad M_3 = \frac{15 + 18 + 22}{3} = 18.3$$

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

$$K_2 = \{6, 8, 9, 12\}$$

$$K_3 = \{15, 18, 22\}$$

∴ Finally we got the final 3 clusters as
from the given dataset :-

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

$$K_2 = \{6, 8, 9, 12\}$$

$$K_3 = \{15, 18, 22\}$$

Group the following objects into $K=2$ groups
of medicine based on the two features
(ph & weight index)

Obj
Medicine A
Medicine B
Medicine C
Medicine D

Attribute 1 (x)

1

2

4

5

Attribute 2 (y)

1

1

3

4

Let's take the two centroid c_1 and c_2
as $(1, 1)$ & $(2, 1)$ respectively. The given
dataset need to be cluster into 2 groups

By Applying the Euclidean distance:

$$D(i,j) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

D	A	B	C	D	Cluster Centroid
0	1	3.6	5	$C_1(1,1)$ - GP1	
1	0	2.82	4.24	$C_2(2,1)$ - GP1	

For C1 (1,1)

$$A = \sqrt{(1-1)^2 + (1-1)^2} \Rightarrow 0$$

$$B = \sqrt{(1-2)^2 + (1-1)^2} \Rightarrow 1$$

$$C = \sqrt{(1-4)^2 + (0-3)^2} = \sqrt{9+4} = \sqrt{13} = 3.6$$

$$D = \sqrt{(1-5)^2 + (1-4)^2} = \sqrt{16+9} = \sqrt{25} = 5$$

For C2 (2,1)

$$A = \sqrt{(2-1)^2 + (1-1)^2} \Rightarrow 1$$

$$B = \sqrt{(2-2)^2 + (1-1)^2} \Rightarrow 0$$

$$C = \sqrt{(2-4)^2 + (1-3)^2} \Rightarrow \sqrt{4+4} \Rightarrow \sqrt{8} \Rightarrow 2.82$$

$$D = \sqrt{(2-5)^2 + (1-4)^2} = \sqrt{9+9} \Rightarrow \sqrt{18} \Rightarrow 4.24$$

The new cluster becomes after finding the distance

$$K_1 = \{A\}$$

$$K_2 = \{B, C, D\}$$

Now, the new centroid as per the new cluster are:

$$C_1(1,1)$$

$$C_2 \left(\frac{2+4+5}{3}, \frac{1+3+4}{3} \right) = C_2(3.67, 2.67)$$

calculating the new euclidean centroid :- distance as per

$$D^* = A$$

	A	B	C	D
O	3.15	1	3.6	5
1	1.36		0.47	1.88

Cluster Centroid

$$C_1(1, 1)$$

$$C_2(3.67, 2.67)$$

The new cluster becomes after finding the
distance :

$$K_1 = \{A, B\}$$

$$K_2 = \{C, D\}$$

Now, the new centroid as per the new cluster are:

$$C_1 = \left(\frac{1+2}{2}, \frac{1+1}{2} \right) = C_1(1.5, 1)$$

ISO 9001:2015 Certified

$$C_2 = \left(\frac{4+5}{2}, \frac{3+4}{2} \right) \Rightarrow C_2(4.5, 3.5)$$

calculating the euclidean distance as per the new centroids.

$$C_1(1.5, 1) \text{ & } C_2(4.5, 3.5)$$

	A	B	C	D	Cluster Centroid
O	0.5	0.5	3.2	4.6	$C_1(1.5, 1)$
1	4.3	3.5	0.7	0.7	$C_2(4.5, 3.5)$

The new cluster becomes after finding the distance.

$$K_1 = \{A, B\}$$

$$K_2 = \{C, D\}$$

Here we can observe the same cluster got resulted as per the previous iteration.

∴ The final clusters group of medicine are

$$K_1 = \{A, B\} \quad \& \quad K_2 = \{C, D\}$$

STUDENT - IDENTIFICATION - TICKET	
Date of Birth:	10/07/98
Date of Birth:	10/07/98
Date of Birth:	10/07/98
Roll No.:	18
Remarks:	AV
Signature of Faculty:	

Subject :- DMBI

(1) (a)

Unknown sample

 $x^* = (\text{age} = \leq 30; \text{Income} = \text{"median"}, \text{student} = \text{"yes"},$
 $\text{credit rating} = \text{"fair"})$

age	Income	Student	credit rating	buys - computer
≤ 30	High	No	Fair	No
≤ 30	High	No	Excellent	No
$31 \dots 40$	High	No	Fair	Yes
≥ 40	Medium	No	Fair	Yes
≥ 40	Low	Yes	Fair	Yes
≥ 40	Low	Yes	Excellent	No
$31 \dots 40$	Low	Yes	Excellent	Yes
≤ 30	Medium	No	Fair	No
≤ 30	Low	Yes	Fair	Yes
> 40	Medium	Yes	Excellent	Yes
≤ 30	Medium	Yes	Excellent	Yes
$31 \dots 40$	Medium	No	Excellent	Yes
$31 \dots 40$	High	Yes	Fair	Yes
≥ 40	Medium	No	Excellent	No

ISO 9001:2015 Certified
NBA and NAC Accredited

STEP1 → LEARNING PHASE .

 $(\text{Buys - computer} = \text{Yes}) = 9/14$ $(\text{Buys - computer} = \text{No}) = 5/14$

age	(buys_computer = Yes)	buys_computer = No
<= 30	2/9	3/5
31..40	4/9	0/5
> 40	3/9	2/5

Income	Yes	No
High	2/9	2/5
medium	4/9	2/5
Low	3/9	1/5

Student	Yes	No
no	3/9	4/5
yes	6/9	1/5

audit rating	Yes	No
Excellent	3/9	3/5
Fair	6/9	2/5

* Given a new instance

x' (age = ' ≤ 30 ', Income = "medium", student = "yes", credit rating = "Fair")

Map Rule.

$$P(x'/y_{\text{es}}) = P(\text{age} = '\leq 30' / \text{yes}) * P(\text{Income} = \text{medium} / \text{yes}) * P(\text{student} = \text{"yes"} / \text{yes}) * P(\text{credit rating} = \text{"Fair"} / \text{yes})$$

Subject :-

 $P(x' | \text{No})$

$$= P(\text{age} = "c = 30" | \text{No}) * P(\text{Income} = \text{medium} | \text{No})$$

$$* P(\text{Student} = \text{Yes} | \text{No}) * P(\text{Credit-rating} = \text{fair} | \text{No})$$

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

$$= \left(\frac{3}{5}\right) \left(\frac{2}{5}\right) \left(\frac{1}{5}\right) \left(\frac{2}{5}\right) \left(\frac{5}{14}\right)$$

$$= 0.0068$$

$$\approx 0.007$$

R

Hence,
we got

$$P(x' | \text{Yes}) = 0.0282$$

$$P(x' | \text{No}) = 0.007$$

Since; $0.028 > 0.007$ ISO 9001:2015 Certified

\therefore The Naive Bayesian classifier predicts buys computer = "yes" for sample x' .