

Date : _____

Name: Mohit sunil surve

Subject: DMBI

RONNO: 74

CLASS : TE/IT

D.O.P	D.O.A	Sign	Marks

Q1. A database has ten transactions. Let minimum support = 30% and minimum confidence = 70%. i) Find all frequent patterns using Apriori algorithm
ii) list strong association rules

(ii) list strong association rules

Transaction id	Items
01	A,B,C,D
02	A,B,C,D,E,G
03	A,B,C,G,H,K
04	B,C,D,E,K
05	D,E,F,H,L
06	A,B,C,D,I
07	B,I,E,K,L
08	A,B,D,E,K
09	A,E,F,H,L
10	B,C,D,F

Soln. S1: Generate single item list set :

$$\text{Support} = \frac{30}{100} \times 100 = 30$$

Items	Support	Item set	above 30% support
A	6	A	6
B	7	B	7
C	6	C	6
D	7	D	7
E	6	E	6
F	3	F	3
G	2		
H	3		
I	1	G	3
J	4	K	4
K	4	L	4
L			

S2: Generate 2 item set:

Item	Support	Item	Support
AB	4	CH	1
AC	4	CK	2
AD	4	CL	1
AE	3	DE	4
AF	1	DF	2
AH	2	DK	1
AK	2	OL	2
AL	2	EF	2
BC	5	EH	3
BD	6	EK	3
BE	4	EL	3
BF	1	FH	3
BH	0	FK	2
BK	3	FL	0
BL	2	HK	2
CD	5	HL	1
CE	2	HL	2
CF	1		1

Item set above 30%. Support

AB	4
AC	4
AD	4
AE	3
BC	5
BD	6
BE	4
BK	3
CD	5
DE	4
EK	3
EK	3

Date : _____

33: Generate 3 item set

Item set of 3 times:

Item Set	Support	Item Set	Support
ABC	3	BEK	2
ABD	4	BEL	1
ABE	2	COE	2
ABK	1	DEK	2
ACD	3	DEL	1
ACE	1		

Item set above 30% support

		Item Set	Support
AEK	1		
AEL	1	ABC	3
ACD	5	ABD	4
BCE	2	ACD	3
BCK	?	BCD	5
BOK	3	BDE	3
BOK	2		

S4: Generate 4 item set

Item set	Support
ABCD	3
ABDE	2
B CDE	2

S6, ABCD is the large item set with minimum support 30%.
following rules are generated :

Rule	Confidence	Confidence %
$A \rightarrow ABCD$	$3/6 = 0.5$	50%

$B \rightarrow ACD$	$3/7 = 0.43$	43%
$C \rightarrow ABD$	$3/6 = 0.5$	50%
$D \rightarrow ABC$	$3/7 = 0.4$	43%
$AB \rightarrow CD$	$3/4 = 0.75$	75%
$BC \rightarrow AD$	$3/5 = 0.6$	60%
$CD \rightarrow AB$	$3/5 = 0.6$	60%
$AC \rightarrow BD$	$3/4 = 0.75$	75%
$AD \rightarrow BC$	$3/4 = 0.75$	75%
$BCD \rightarrow A$	$3/5 = 0.6$	60%
$ACD \rightarrow B$	$3/3 = 1$	100%
$ABD \rightarrow C$	$3/4 = 0.75$	75%
$ABC \rightarrow D$	$3/3 = 1$	100%

From above rules generated, only the rules having greater than 101. are considered as final rules as so, final rules are

$$AB \rightarrow CD$$

$$AC \rightarrow BD$$

$AD \rightarrow BC$

$A \rightarrow B$

ABO → C

$$\underline{ABC \rightarrow D}$$

Q2. Use Apriori Algorithm with mini-support count = 2, confidence = 60% to find all frequent itemsets and strong association rules. Consider the transaction database given below

Transaction Id	Items
10	1,3,4
20	2,3,5
30	1,2,3,5
40	2,5
50	1,3,5

Soln SI: Generate single item set

Item set	Support	Item set above	2 support count
1	3	1	3
2	3	2	3
3	4	3	4
4	1 9 6 6 6 6 6	5	4
5	4		

S2: Generate 2 time set

Item Set	Support
1, 2	1
1, 3	3
1, 4	1
1, 5	2
2, 3	2
2, 5	3
3, 4	1
3, 5	3

item-set above \rightarrow support count

item set Support

Date :

1,3	3
1,5	2
2,3	2
2,5	3
3,5	3

S3: Generate 3-item set

item-set	Support
1, 2, 3	1
1, 2, 5	1
2, 3, 5	2
1, 3, 4	1
1, 3, 5	2

$\{1, 3, 5\}$, $\{2, 3, 5\}$ is large item-set with minimum support 12'.

Following rules are generated:

Rule	Confidence	Confidence %
$2^1 3 = 5$	$2 2 = 1$	160%
$3^5 = 2$	$2 3 = 0.66$	66%
$2^1 5 = 3$	$2 3 = 0.66$	66%
$2^1 3^1 5$	$2 3 = 0.66$	66%
$3 = 2^1 5$	$2 4 = 0.5$	50%
$3 = 2^1 3$	$2 4 = 0.5$	50%

from above generated rules, only rules having sol. are considered as final rules.

\therefore final rules are

$$2^1 \cdot 3 = 5$$

$$3^1 \cdot 5 = 2$$

$$2^15 = 3$$

$$2 = 3^{\wedge}5$$

Rules generated

Rules	Confidence	Confidence
$1 \wedge 3 = 5$	$2/3 = 0.66$	66%
$3 \wedge 5 = 1$	$2/3 = 0.66$	66%
$1 \wedge 5 = 3$	$2/2 = 1$	100%
$1 = 3 \wedge 5$	$2/3 = 0.66$	66%
$3 = 1 \wedge 5$	$2/4 = 0.5$	50%
$5 = 1 \wedge 3$	$2/4 = 0.5$	50%

Q3 Apply Naive Bayes classifier algo, to classify an unknown sample x (outlook = sunny, Temp = cool, Humidity = High, windy = false).

Sample dataset is given as follows:-

outlook	Temp	Humidity	windy	Class
Sunny	Hot	high	false	N
Sunny	Hot	high	true	N
Overcast	Hot	high	false	P
Rain	mild	high	false	P
Rain	cool	normal	false	P
Rain	cool	normal	true	N

Overcast	cool	normal	true	P	
Sunny	mild	high	true	N	
Sunny	cool	normal	false	P	
Rain	mild	normal	false	P	
Sunny	mild	normal	false	P	
Overcast	mild	high	true	P	
Overcast	hot	normal	true	P	
Rain	mild	high	false	N	

Soln Given a training set we can compute the probabilities as follows:

- The classification prior may be formalized using a posterior probability of the classes: $p(c_i | \mathbf{x})$ is probability of the i^{th} class given sample tuple $\mathbf{x} = (x_1, x_2, \dots, x_n)$ is of C_i .

from the above given sample data, calculate probability for play tennis (P) and tennis (n) for all attribute.

OUTLOOK

$P(\text{sunny} \text{Yes}) = 2/9$	$P(\text{sunny} \text{No}) = 3/5$	
$P(\text{overcast} \text{Yes}) = 4/9$	$P(\text{overcast} \text{No}) = 0$	$P(\text{Yes}) = 9/14$
$P(\text{rain} \text{Yes}) = 3/9$	$P(\text{rain} \text{No}) = 2/5$	$P(\text{No}) = 5/14$

Temp

$$P(\text{hot}|\text{yes}) = 2/9 \quad P(\text{hot}|\text{no}) = 2/5$$

$$P(\text{mild} | \text{yes}) = 419 \quad P(\text{mild} | \text{no}) = 215$$

$$P(\text{cold} | \text{Yes}) = 3/9 \quad P(\text{cold} | \text{No}) = 1/15$$

Humidity

$$P(\text{high} \mid \text{YES}) = 3/9 \quad P(\text{high} \mid \text{NO}) = 4/15$$

$$P(\text{False} | \text{Yes}) = 619 \quad P(\text{False} | \text{No}) = 215$$

An unseen sample $y = \text{strain, not, high, false}$

$$\begin{aligned}
 P(Y|Yes, P(Yes)) &= P(\text{rain}|\text{Yes}) \cdot p(\text{not } Y|\text{Yes}) \cdot P(\text{high}|\text{Yes}) \\
 &\quad \cdot P(\text{false}|\text{Yes}) \cdot P(\text{No}) \\
 &= 3/9 \cdot 2/9 \cdot 3/9 \cdot 6/9 \cdot 9/14 \\
 &= 0.010582
 \end{aligned}$$

choose the class so that it maximises the probability. This means new instance can be classified as no sample is classified as class No.

An unsee sample = <sunny, cool, high, true>

$$\begin{aligned}
 P(Y|yes) \cdot P(yes) &= P(\text{sunny}|yes) \cdot P(\text{cool}|yes) \cdot P(\text{high}|yes) \\
 &\quad \cdot P(\text{true}|yes) \cdot P(yes) \\
 &= 219,319 \cdot 319 \cdot 9114 \\
 &\approx 0.0053.
 \end{aligned}$$

$$P(\text{NO}) \cdot P(\text{y/NO}) = P(\text{sumy/NO}) \cdot P(\text{cool/NO}) \cdot P(\text{high/NO}) \cdot P(\text{true/NO}) \\ \cdot P(\text{NO}) \\ = 3/5 \cdot 1/5 \cdot 4/5 \cdot 3/5 \cdot 5/4 \\ = 0.0206$$

Now, choose the class so that it maximize this probability. This means that the new instance will be classified as no.

Q4. Define classification discuss the issues in classification. A simple ex. from stock market involving only discrete ranges has profit as categorical attribute, with value {up down} and the data is

Age	competition	Type	Profit
old	yes	software	down
old	no	software	down
old	no	software	down

		mid	yes	Software	Down
		mid	yes	hardware	Down
		mid	No	hardware	Down
		mid	No	Software	Up
		new	yes	Software	Up
		new	No	hardware	Up
		new	No	Software	Up

Soln classification

- it is the form of data analysis
 - classification constructs classification model based on training data set
 - using this model it classifies the new data.
 - classification model it classifies the predicate.
 - categorical class labels
 - for e.g. we can classify model to categories bank loan application as either safe or risky.

Issues in classification

- ① Data cleaning - removing the noise.
 - ② Feature Relevance analysis - Also & relevant attributes
 - ③ Data transformation & reduction.
 - a) Normalization - involves scaling all values.
 - b) Generalization - transformed by generation.

class = Profit p=up , n=down

Total No. of records = 10

$$\therefore p = 5 \text{ & } n = 5$$

$$\text{so info.gain} = I(P|N) = -\frac{P}{P+N} \log_2 \frac{P}{P+N} - \frac{N}{P+N} \log_2 \frac{N}{P+N}$$

$$\therefore I(5,5) = \frac{-5}{10} \log_2 \frac{5}{10} - \frac{5}{10} \log_2 \frac{5}{10}$$

$$I(5,5) = 1$$

S1. Compute entropy for age : cold, mild, new) for age = old.

pi with 'up' class = 0 & n, with 'down' class = 3

$$I(P_1, P_2) = I(0, 3) = -\frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right)$$

Similarly for different age

ranges $I(p_i, n_i)$ are calculated as follows:

Age	p_i	n_i	$I(p_i, n_i)$
old	0	3	0
mid	2	2	1
new	3	0	0

$$E(A) = \sum_{i=1}^v p_i + n_i I(p_i, n_i)$$

$$E(\text{Age}) = \frac{3(0)}{10} + \frac{4(1)}{10} + \frac{3(0)}{10} = 0.4$$

$$\text{Gain}(I, \text{age}) = I(\rho, n) - E(\text{age}) \\ = 1 - 0.4 = 0.6$$

Just as we calculated entropy & gain for age, same can be calculated for competition & type as follows

Entropy can be competition (yes, no)

competition	can be compet'	p _i	n _i	I(p _i n _i)
yes		1	3	0.811
no		4	2	0.918

$$E(\text{competition}) = \frac{4}{10} \times (0.811) + \frac{6}{10} (0.918) = 0.875$$

$$\text{Gain (1-competition)} = 1 - 0.875 \\ = 0.125$$

Entropy of type : $(S|\omega, H|\omega)$:

Type	p_i	n_i	$I(p_i, n_i)$
S/ ω	3	3	1
H/ ω	2	2	1

$$E(\text{type}) = \frac{6}{10}(1) + \frac{4}{10} \times 1 = 1$$

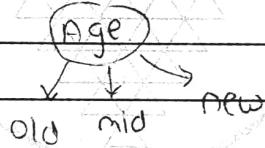
$$\therefore G(T, \text{type}) = 1 - 1 = 0$$

$$\text{Gain}(T, \text{age}) = 0.6$$

$$\text{Gain}(\Gamma, \text{competition}) = 0.125$$

$$\text{Gain}(\tau, \text{type}) = 0$$

∴ we get



S2: Since attribute age is not, we have to decide remaining two attributes for age branch node.

Age	Contest	Type	Profit
old	yes	Slow	down
old	no	Slow	down
old	no	High	down

\therefore The profit is down for every old age is simple branched as down.

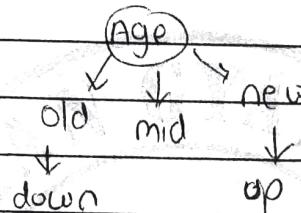
S3: consider age = new

Age competition Type Profit

New	yes	S W	up
new	no	B W	up
NEW	no	S W	up

\therefore profit is up for every new age is simply branched as up.

so, new decision tree becomes



S4: Now consider $\text{age} = \text{mid}$.

Age	competition	Type	Profit
mid	yes	S/w	down
mid	yes	H/w	down
mid	No	H/w	up
mid	No	S/w	up

(i) Compute Entropy for competition : (yes, no)

$$p \approx 1, n \approx 2$$

competition	p_i	n_i	$E(p_i, n_i)$
yes	0	2	0
no	2	0	0

$$\therefore \text{Entropy (competition)} = \frac{2}{4} (0) + \frac{2}{4} (0) = 0$$

$$\text{Gain (age, competition)} = 1 - 0 = 1$$

(ii) compute Entropy for type $(S|w, H|w)$

Type	p_i	n_i	$I(p_i, n_i)$
S1w	1	1	1
H1w	1	1	1

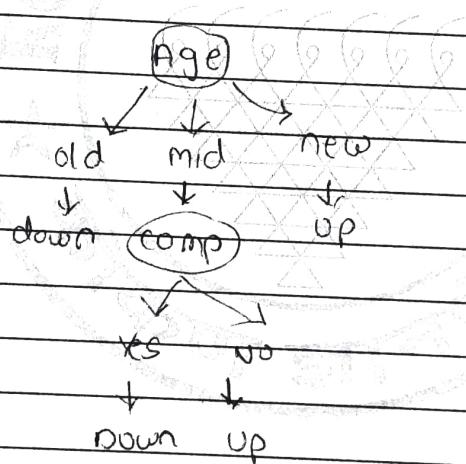
$$E(\text{type}) = \frac{2}{4} x(1) + \frac{2}{4} x(1) = 1$$

$$\therefore G(\text{Fage}, \text{Type}) = 1 - 1 = 0$$

Grain (Tage, comptⁿ) = 1

Gain (Tage, Type) = 0

So, we got decision tree as follows



Q5.

What is clustering? Explain k-means clustering algorithm.
 Suppose data for clustering is $\{2, 4, 10, 12, 03, 20, 11, 5, 25, 36, 41, 14\}$. Assuming number of clusters to be 2 i.e. $k=2$, cluster given data using above algorithm.

|| 80 | n

- Clustering is an unsupervised learning problem.
 - Clustering is a data mining technique used to place data.

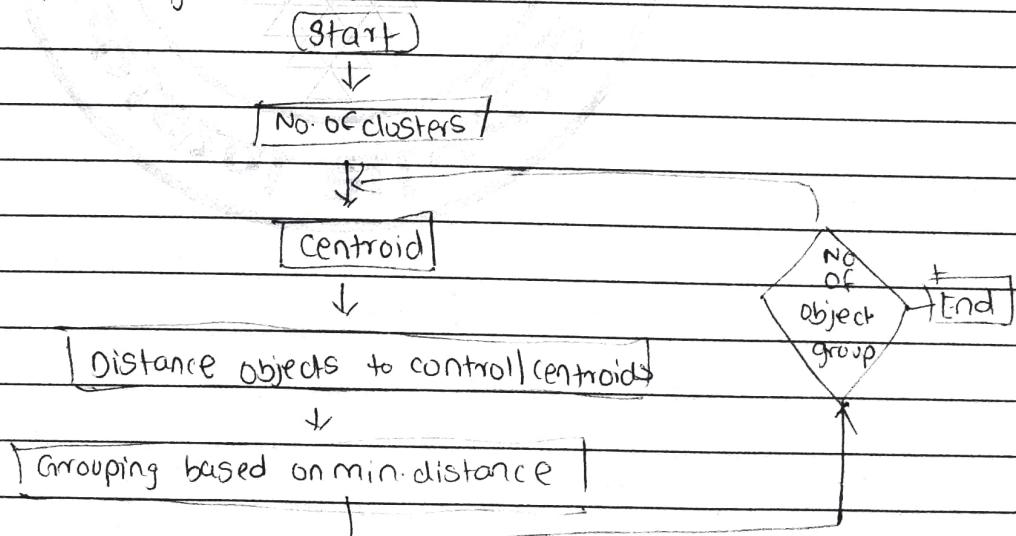
Date : _____

elements into related groups without advanced knowledge of group definition.

- It's a process of partitioning data objects into sub classes which are called as clusters
 - A cluster contains data objects which have as high inter similarity and low intra similarity.

K-means clustering algorithm:

- K-means clustering algorithm is one of the partitioning method. It's simplest unsupervised learning methods.
 - K-means clustering aims to partition n' observations into k' clusters in which each observation belongs to the clusters with the nearest mean, serving as a prototype of the cluster.
 - This results in partitioning of data space.
 - K is the integer number.



Flowchart of K-means

Above figure shows flowchart of K-means clustering

- Define K centroids for K clusters which are generally far only from each other.
 - Then group elements into clusters which are nearer to the centroid of that cluster.
 - After this 1st step, calculate new centroid from each cluster based on elements of that cluster.
 - Follow the same method & group the elements based on new centroid.

Given :- {2, 4, 10, 12, 3, 20, 11, 25, 5, 36, 41, 14} LCM = ?

Randomly assign means: $m_1 = 4$, $m_2 = 12$

Again, calculate new mean for new cluster groups.

$$\text{① } k_1 = \{2, 4, 3, 5\} \quad k_2 = \{10, 12, 20, 11, 25, 36, 41, 14\}$$

$$\textcircled{2} \quad K_1 = \{2, 4, 3, 5, 10, 12, 11\} \quad K_2 = \{20, 25, 36, 41, 14\}$$

$$(3) \quad K_1 = \{24, 3, 5, 10, 12, 11, 14\} \quad K_3 = \{20, 25, 36, 41\}$$

$$\textcircled{4} \quad K_1 = \{2, 4, 3, 5, 10, 12, 11, 14\} \quad K_2 = \{20, 25, 30, 41\}$$

Stop as clusters in step 3 & 4 are same. Clusters in last two groups are identical.

$$\text{so, } k_1 = \{2, 4, 3, 5, 10, 12, 11, 14\} \text{ and } k_2 = \{20, 25, 36, 41\}$$

Q6.

Explain Business Intelligence Issues

→ Organizations and People 1.

- Management within an organization do not agree that the decision taken based on data or evidence work for them, they prefer to run the operation from instinct.

2. In order to access business progress there is no overall business strategy laid out with objectives and measures for those objective.
 3. The data needed for business intelligence system, cannot be obtained from IT personal as they accomplished and they have no resources available.
 4. The performance at the business either mating ace of BI or not there is no incentive provided to the staff within the organization.
 5. There is no obvious time to establish a BI system. The business is in a static at high chance or flux.

2) Data and Technology:

- 1. The data at the organization is not clean the time and effort needed to correct or handle this type of data leads to an unsuccessful delivery at the BI support. For e.g. there could be many different destinations for the same item.
e.g. a customer may be coded differently in sales system than held in accounts system.
 - 2. The technology chosen for BI turns out to be so particular that it ultimately results in time consuming process which leads to delay in project completion.
 - 3. The BI technology discourages the use of system due to the following reasons:
 - The info. presentation quality is poor or limited.
 - The response time for data presentation is slow.
 - It is too difficult for data presentation to cut new question of BI technology for either of them and users of BI expect