

## Nappayya Eadiyappa

CSE181005

## Data Mining

June - 2016

Q1) What is Data Mining? Explain any four data visualization techniques with examples.

### Data Mining

Data mining is the process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems. Data mining is the step of the knowledge discovery in data base process or KDD.

\* Four data visualisation Techniques :-  
Data visualisation is the process of conveying information in a way that can be quickly digested by the viewer. Examples are every where, and see the daily charts, graphs, digital images & movies.

### 2020 Techniques:-

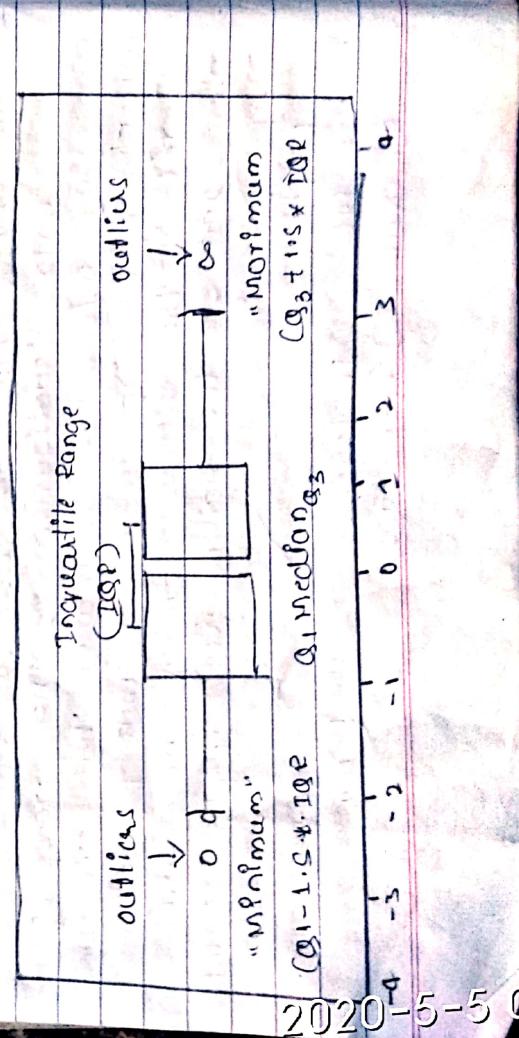
- 1) Box plots.
- 2) Histograms.
- 3) Heat maps.
- 4) Charts.
- 5) Tree maps.

09:32

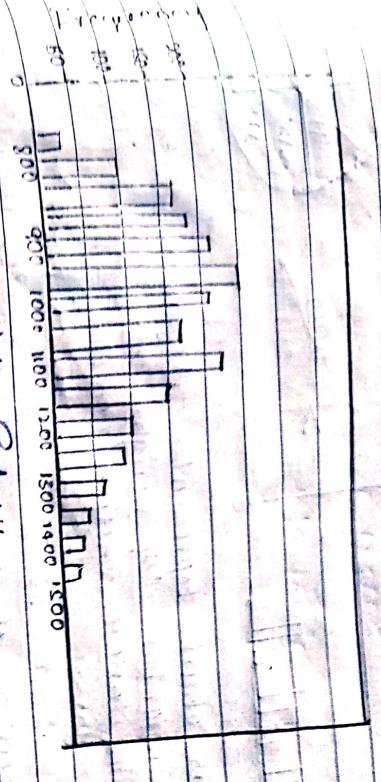
### D Box plots:-

A method for graphically depicting group of numerical data through their quartiles the distribution of data based on 5 features.

- 1] Minimum - minimum value in the dataset
- 2] First quartile - the middle value between the smallest value & the median of the dataset
- 3] Median - the middle value between the max of the dataset
- 4] Third quartile : The middle value between medians of the highest value of the dataset
- 5] Maximum : The maximum value in the dataset also the lower & upper quartile are shown as horizontal line of the rectangle it has vertical line outside to indicate the median value



## Histograms Histograms of Monthly Salary



### gross monthly salary

is an accurate representation of the distribution of numerical data if it relates only one button of includes bin or break the range of variable that is divide the entire range of values that into a series of 10-10 values that count how many values fall in each interval.

3) charts, gives graphical representation of data

It has various types.

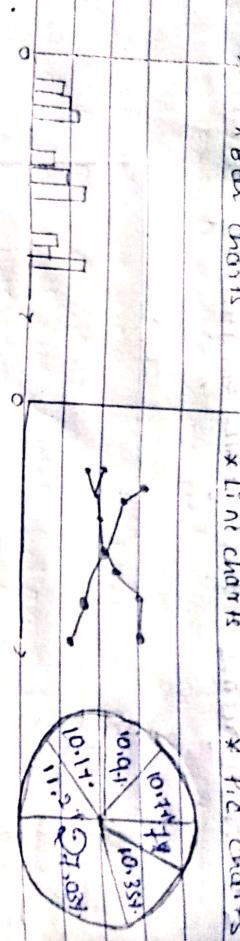
\* Bar charts \* Line charts \* pie charts

2020-5-5 09:39

\* Bar charts

\* Line charts

\* Pie charts



\* Bar chart Represent categorical data with rectangular bars. of height proportional to the values they represent

\* Line chart which displays information as a series of data points called markers connected by straightline segment

\* pie chart -  
is a circular statistical graph  
which divided into slices to illustrate numerical proportions.

- ④ Tree maps It's a method of displaying hierarchical data using nested figures like rectangle has a tree like structure with each branch is given a rectangle which the tiles which smaller rectangle representing sub branches.  
A leaf node rectangle has an area proportional to specified dimension of data

2020-5-5 09:39

How do you handle missing values in a table of different types of attributes?  
The missing value problem for the data set given below

Floor Area	Rental Price	Type
42	10000	1 BHK
52	13000	?
52	11000	2 BHK
52	?	2 BHK
65	18000	2 BHK
?	13500	3 BHK
60	12200	3 BHK

Attribute descriptions: Floor Area is numeric.

Rental Price: Double.

Type: categorical.

→ Missing values. In a dataset of different types of attributes.

- 1) Ignore the data row
- 2) Use a global constant to fill in for missing values
- 3) Use attribute mean
- 4) Use attribute mean for both samples belonging to the same class.
- 5) Use a data mining algorithm to predict the most probable value.

2020-5-5 09:39

Q. A. Briefly compare the following concepts  
i) snowflake schema and fact constellations

snowflake schema	Fact constellations
1) snowflake schema contains the large central fact table dimension tables & sub dimensions.	1) While in fact constellations schema dimension tables are shared by many fact tables.
2) This schema gives significant storage.	2) This schema does not give storage.
3) This schema consists of one fact schema and a time dimension.	3) This schema consists of more than one fact schema at a time.
4) Tables can be maintained easily.	4) Tables are tough to maintain.
5) This schema is norinalized from a fact star schema.	5) This schema is a flat schema for star schema.
6) This is easy to operate as compared to fact constellations.	6) In this schema is not easy to operate as compared to snowflake schema as it has multiple normal joins between the tables.

2020-5-5 09:38

Q.4) How does the ETL in this help to access the data from database.  
Ans:- From database, simple + less complex query is used.

Q.5) Discovery driven cube is used, for guiding the user for data analysis process by making use of the pre-computed data + virtual cubes to indicate exception at cell level aggregation via virtual warehouse.

Q.6) Regarding the computation of measures, in a data cube, we have three categories of measures. With an example for each category, measures. Data cube computation is an essential part in data warehouse implementation!

The cube is used to represent data along some measure of interest although called a "cube" it can be 2-dimensional 3-dimension or higher dimensional

Ex:- We have a data base that contains transaction information relating company sales of a part to a customer at a store location

to enumerate their categories of measures based.

2020-5-5 09:38

On the kind of aggregate functions used in computation of data cube they are -

- 1) Distributive
- 2) Algebraic
- 3) Holistic

Distributive :-

It is the result derived by applying the function to an aggregate by applying the function to  $n$  values which is the same as that derived by applying the function on all the data without partitioning.

Examples for count(), sum(), min(), max(),

2) Algebraic :-

If it can be computed by an algebraic function with  $m$  arguments. (Where  $m$  is a bounded integer) each of which is obtained by applying a distributive aggregate function.

Examples - avg(), min(), standard-deviations,

3) Holistic :-

It is the result derived by applying the function to all the data in the storage size needed to describe a subcube.

Examples -

Medians, modes, rank(),

diff. types to the frequent item set mining.  
With the following rule 10th on example.

### Support

If  $\mu$  is one of the measure of interestingness  
like rule about usefulness of transaction  
of rules. The support measure total. 54 of  
transaction in database follow the rule.  
 $\text{Support}(X \rightarrow Y) = \text{Support} - \text{count}(A \cup B)$

### Confidence

A confidence of 60% means that 60%  
of the customers who purchased A in fact and  
bought also bought B.

### Confidence $(A \rightarrow B)$ Support - count(A)/Support - count(B)

If a rule's confidence both minimum support  
and maximum confidence by far a strong  
rule.

### Example:

Support - count(X) / number of transactions.  
Set where X appears. If X is a union B  
then it is the number of transactions in  
which A and B both are present.

### Frequent Itemset

the frequent item-set mining. is an interesting  
branch of data mining that focuses on looking  
at sequences of actions or events that base  
data tokens that occurs at sets of instance. for

that each has a number of features.

A typical example of frequent itemset mining is market basket analysis.

Example: Two item sets.

T<sub>1</sub> {A,B,C,D}

T<sub>2</sub> {A,D}

T<sub>3</sub> {A,E}

T<sub>4</sub> {C,E}

Now if the minimum support threshold is 50% i.e. it is present in 3 so. 75%. C, D, E are present in 2 transactions. Satisfying 50%. So they are frequent.

v) closed itemsets-

It is a frequent itemset for which none of its immediate supersets have that same support count as itself.

Example of above, illustration.

{A,B} is closed because none of its

Superset have that same support as. i.e. {A,B,C}, {A,B,D}, {A,B,C,D} are closed. frequent itemset {A,B,C,D} is not closed because its immediate superset {A,B,D} also has the same support count 2.

v) Association Rule:

An association rule is pattern

that states, when an event occurs another event occurs with certain probability

Association rules: find all sets of items

2020-5-5 09:38

that have support greater than the minimum support than using the long items to generate the desired rules that have confidence greater than the minimum confidence.

### v) Frequent Pattern (FP)

Frequent pattern tree is a tree-like structure that is made with the initial item sets of the database. The purpose of the FP tree is to mine the most frequent pattern. Each Node of the FP tree represents an item of the item sets.

FP tree - compress data base onto tree.

I.C

Ex T.D Item Sets.  
T<sub>1</sub> {M,O,N,E,Y}   
T<sub>2</sub> {D,O,N,K,E,Y}   
T<sub>3</sub> {M,A,K,E,Y}   
T<sub>4</sub> {M,U,C,K,E,Y}   
T<sub>5</sub> {C,O,O,K,I,E,Y}

Item	Support count	L
M	3	E
O	3	C
N	2	M
R	5.	3
E	4	0
Y	3	4
D	1	3
A	1	2
V	1	1
C	2	0
I	1	1

frequent pattern.

3.b

ordered items.

Item ID	S-C	N.L
E	5	- - - - E 84 - 5
E	4	- - - - E 62 - 4
M	3	- - - - M 52 - 3
O	3	- - - - O 81 - 3
Y	3	- - - - Y 61 - 3

2020-5-5 09:38

FPTree	conditional.	Frequent pattern
FP tree.	4 < E, Y & 3 &	
{E, 3}	0 < R, 0 ; 3 > < E, 0 ; 3 >	
{E, 3}	< 0, C, 1 & 3 >	
{E, 3}	M < M, 1 & ; 3 >	
{E, 3}	E < E, R ; 3 >	

- Q) How do you construct FP-tree effectively  
to use in the FP-growth algorithm?  
Discuss with an example

### FP Growth Algorithm:

This algorithm is an improvement to the Apriori method. A frequent pattern is generated throughout the need for candidate generation -ed FP Growth Algorithm represents the database in the form of a tree called frequent pattern tree or FP tree

This tree structure maintains the association between the itemsets, the database e.g. frequent using one frequent item tree fragment part as called a pattern fragment each node of the FP tree represents an item of the itemsets. The Root Node represents an item of the itemsets. The Root Node represents null while the lower nodes represents the item sets. The association of the node with the lower nodes that is the item sets with the item sets are maintained while forming

Table 2.

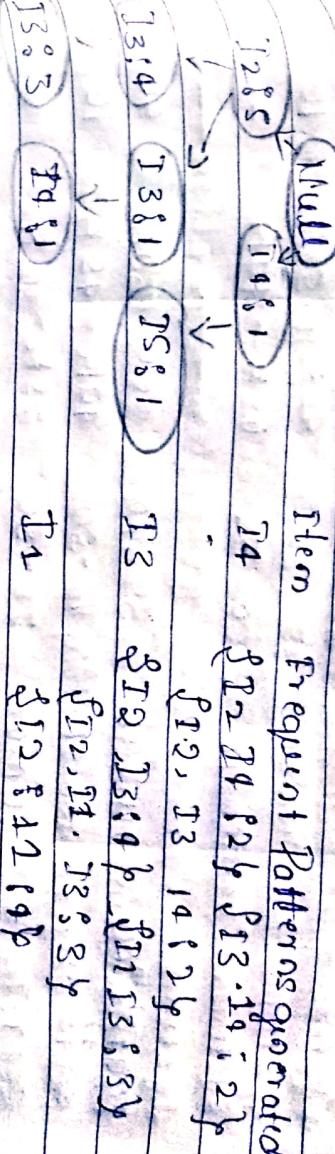
Item	count
D <sub>1</sub>	4
D <sub>2</sub>	5
D <sub>3</sub>	2
D <sub>4</sub>	9
D <sub>5</sub>	2

Q) Sort the longest insuring order

Table - 3  
Item count      Item conditional path length      condition

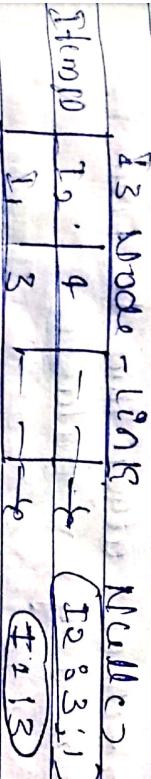
Item	count	Item	conditional path length	condition
T <sub>1</sub>	5	T <sub>4</sub>	1	[P <sub>2</sub> P <sub>1</sub> → S <sub>1</sub> :1]
T <sub>2</sub>	5	T <sub>5</sub>	1	[P <sub>2</sub> P <sub>2</sub> → S <sub>1</sub> :2]
T <sub>3</sub>	5	T <sub>6</sub>	1	[P <sub>2</sub> P <sub>3</sub> → S <sub>1</sub> :3]

### Build FP-tree



$T4:1$

The diagram given below depicts the conditional FP-tree associated with the conditional node



4 c) The following table consists of training data from an employe database. The data have been generalized. For example "30-35" for age is represented - the age range of 31 to 35 for a given row entry. count represents

the No. of data tuples having the values for department, status, age & salary & given in Row

Department	Status	Age	Salary	Count
Sales	Senior	31-35	46K-50K	30
Sales	Junior	26-30	35K-40K	90

Systems	Senior	35 - 40	50K - 65K	15
systems	Junior	30 - 35	45K - 50K	10
systems	Senior	30 - 35	55K - 70K	20
systems	Junior	35 - 40	55K - 70K	—
Marketing	Junior	35 - 40	50K - 65K	40
secondary	Senior	45 - 50	40K - 45K	10
secondary	Junior	25 - 30	35K - 50K	15
(i)	Senior	45 - 50	45K - 50K	3

Let  $\hat{S}_k$  be the class label constraint  
of design tree from the given data

to use K-means algorithm to form the  
clusters from a given dataset

=> K-means algorithm is a iterative algorithm  
that tries to partition the dataset into K  
distinct non overlapping sub groups

to only group

The way K-means algorithm works is as  
follows:-

Step 1 -> specify number of clusters K

Step 2 -> initialize centroids by first stuff

the dataset and then randomly selecting  
K data points for the centroids without replace-

ment

Step 3 -> keep iterating until there is no

change to the centroid assignment  
of data points to clusters until changes

2020-5-5 09:38

compute the sum of the squared distance between data points and all centroids assign each data point to the closest cluster centroid

compute the centroid for the clusters by taking the average of the all data points that belong to each cluster.

The approach known as follows to solve the problem is called expectation maximization

This algorithm alone at minimizing the objective function given by error function given by

$$J(C) = \sum_{i=1}^C \sum_{j=1}^{n_i} (c_j - x_i)^2$$

where  $\|x_i - c_j\|^2$  is the euclidean distance between  $x_i$  &  $c_j$   $c_j$  is the mean of points in  $j$ th cluster  $c_j$  is the mean of cluster unless

- 5 a.) What are outliers? Address the various challenges of outlier detection process  
An outlier is an object that deviates significantly from the rest of the objects. They can be caused by measurement or execution errors to as outlier analysis can outliers mining  
Outliers ≠ Noise data.

2020-5-5 09:38

- 4 challenges of outliers detection
- 1) Modeling normal objects and outliers properly
  - 2) Application-specific outlier detection
  - 3) Handling noise in outlier detection
  - 4) Understandability

2020-5-5 09:38

- 1) modeling normal objects & outliers properly!
- Hard to enumerate all possible Normal behaviours in an application. The border between normal & outlier objects is often a gray zone.

- 2) Application-specific outlier detection
- choice of distance measure among objects  
and the model of relationship among objects  
are application dependent

For clinic data a small deviation could be an outlier while in marketing analysis large fluctuations

- 3) Handling noise in outlier detection
- noise is ~~disturbing~~ disturbing about 40%  
noise many distort the normal objects & blur the distinction the normal objects & outliers  
it may help hide outliers & reduce effective of outlier detection.

- 4) Understandability

Understandability why these are outlier

and outliers :- the unlikelihood of the object being generated by a Normal mechanism.

Suppose a city's average temperature (in Celsius) in July in the last 10 years are

$$28.0, 28.9, 28.9, 29.0, 29.1, 29.1, 29.2$$

$$29.2, 29.3, 29.4$$

Identify outliers using Maximum Likelihood method.

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n-1} x_i^2$$

$$\hat{x} = 29.0 + 28.9 + 28.9 + 29.0 + 29.1 + 29.1 + 29.2 + 29.2 + 29.3$$

$$+ 29.4$$

$$= 10$$

$$= 28.61$$

$$\hat{\sigma}^2 = (29.0 - 28.61)^2 + (28.9 - 28.61)^2 + (28.9 - 28.61)^2 +$$

$$(29.0 - 28.61)^2 + (29.1 - 28.61)^2 + (29.1 - 28.61)^2 +$$

$$(29.2 - 28.61)^2 + (29.2 - 28.61)^2 + (29.3 - 28.61)^2$$

$$+ 10$$

$$= 21.25 + 0.0891 + 0.0891 + 0.1521 + 0.2461 + 0.2461 +$$

$$0.3481 + 0.3481 + 0.4761 + 0.6241$$

$$10$$

$$= 23.8469$$

$$\hat{\sigma}^2 = \sqrt{2.3} = \sigma^2 = 1.811$$

6(a) Briefly outline how to compute the dissimilarity between object described by the following.

i) Nominal attributes

A categorical variable is a generalization of the binary variable so that it can have one or more than two states. The dissimilarity between two objects can be computed based on the ratio of mismatches;  $d_{ij} = p - \rho_{ij}$

where  $m$  is the no. of matches

Give the no. of variables for which  $j$  and  $k$  are in the same state &  $p$  is the total no. of variables

Ex:	Object	test-1	test-2	test-3
Telephones	Nominal	ordinal	numerical	
1	Code A	Excellent	4.5	
2	Code B	Tall	2.2	
3	Code C	Good	6.4	
4	Code D	Excellent	2.8	

only object-identifier & the attribute test-1 are available when test-2 is ~~the~~ Nominal

0	$d(2,1)$	0	
$d(3,1)$	$d(3,2)$	0	
$d(4,1)$	$d(4,2)$	$d(4,3)$	0
0	.	.	

1	0	
0	1	1

2020-5-5 09:38

From this we see that these objects

are dissimilar enough objects  $1 \& 4$  (i.e.)

$$d(1, 4) = 6$$

iii) Asymmetric binary attributes.

A binary attribute is asymmetric if one of the states are not equally important such as the positive & negative outcome of a medical test for HIV test (say) the most important outcome (likely case) is, whereas the next one is '0'

ii) If all binary variables have the same weight we have the contingency Table

object

	1	0	sum
1	1	2	3
0	0	3	3
sum	1	5	6

To computing the dissimilarity between any two state binary variables the role of negative may be considered unimportant thus it is ignored in the computation that is

iii) Numeric attributes -

use euclidean distance manhattan distance or supremum distance Euclidean distance is defined as

The Manhattan City block distance is defined as

$$d(x, y) = |x_1 - y_1| + |x_2 - y_2| + \dots + |x_n - y_n|$$

### iv) Term Frequency Vectors -

To measure the distance between complex objects represented by vectors it is often easier to abandon traditional metric distance computation & introduce a non-metric similarity function. The similarity between two vectors  $x$  &  $y$  can be defined as cosine measure as follows

$$S(x, y) = \frac{x^T y}{\|x\| \|y\|}$$

Where  $x^T$  is a transposition of vector  $x$ .  
 $\|x\|$  is the Euclidean norm of vector  $y$  i.e. is the essentially the cosine of the angle between vectors  $x$  &  $y$ .

### Q b) Explain data warehouse design process with example

A data warehouse is a single data repository where a record from multiple data sources is integrated for online business analytical processing (OLAP). This implies a data warehouse needs to meet the requirement from all that business business stages within the entire organization. Thus data warehouse design is highly complex, lengthy & never easy or simple process.

Data warehouse design takes a method different from new materialization: In the industries it sees data warehouses as data barn

framed for need (cr) to be organized in  
a data base as the data ware house

There are two approaches

- 1) "Top-down" approach
- 2) "bottom up" approach

D Top down Design approach

To this design a data warehouse is describable  
as a subject-oriented time variant non-volatile  
& integrated data repository for the entire enterprise.  
data from different sources are validated  
performed saved in a normalized (up to 3NF)  
data base as the data ware house.

advantages :-

D data marts are loaded from the data warehouses  
very easy

discarding e-

& the data masters are denoted on the bottom-up approach design

4(a) A database has 4 inv + transactions. let Inv =

Scpport = 60.1. & Min - conf = 80, TID

T100, T200, T300

Y M A E Y

T400, S M V C K H Y

T500, E C O O F D E Y

Find all frequent item sets using Apriori

Implementation

C1

C2

Item count

Item count

Set 1

Set 2

Set 3

Set 4

Set 5

Set 6

Set 7

Set 8

Set 9

Set 10

Set 11

Set 12

Set 13

Set 14

Set 15

Set 16

Set 17

Set 18

Set 19

Set 20

Set 21

Set 22

Set 23

Set 24

Set 25

Set 26

Set 27

Set 28

Set 29

Set 30

Set 31

Set 32

Set 33

Set 34

Set 35

Set 36

Set 37

Set 38

C3

Item set count

EHO 3

AMO 7

AMY 2

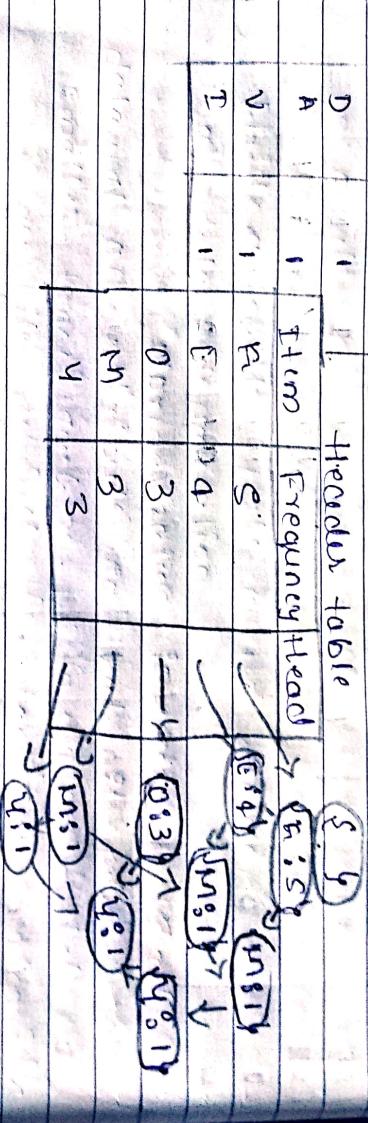
HOY 2

### FP Growth

09:37

Item	Frequency	Item	Frequency	Item	Frequency
B	4	T100	1	MONEY	1
C	3	T200	1	D.O.N.K.E.Y	1
E	3	T300	1	M,A,B,E	1
F	2	T400	1	MUCKY	1
G	2	T500	1	COKE	1
H	1				
I	1				
J	1				
K	1				
L	1				
M	3				
N	3				
O	3				
P	3				
Q	3				
R	3				
S	3				
T	4				
U	4				
V	4				
W	3				
X	3				
Y	3				
Z	3				

Header Table



Item conditional frequent itemset frequent pattern

Pattern base: the pattern base generated

E frequent itemset C.R.E.P.

O frequent itemset C.R.O.P.

M frequent itemset C.M.P.

Y frequent itemset C.R.Y.P.

(K,E,O,I) frequent itemset C.K.E.O.I.P.

(E,O,I,D) frequent itemset C.E.O.I.D.P.

(R,M,I,D) frequent itemset C.R.M.I.D.P.

list all the association rules matching the following metarule. Where x is a variable representing customer & item denotes variable representing items (ex - B, B etc) for all x transaction

## Performance

Metric	Formula	Evaluating focus
Accuracy (Acc)	$\frac{tp + tn}{tp + fn + fp + tn}$	In general the accuracy metric measures the ratio of correct predictions over the total number of instances evaluated.
Error Ratio	$\frac{fp + fn}{tp + fp + fn + tn}$	Missclassification error measures the ratio of incorrect predictions over the total number of instance evaluated.
Sensitivity	$\frac{tp}{tp + fn}$	This metric is used to measure the fraction of negative patterns that were correctly classified.
Precision (P)	$\frac{tp}{tp + fp}$	This metric is used to measure the fraction of negative patterns that were correctly classified.
Recall (R)	$\frac{tp}{tp + fn}$	Recall is used to measure the fraction of true patterns that were correctly classified.
F - Measure	$\frac{2 * P * R}{P + R}$	This metric represents the harmonic mean between Recall and precision values.
Geometric Mean	$\sqrt{P * R}$	The metric is used to maximize the product of true positive rate and recall simultaneously.
Mean GM	$\frac{P + R}{2}$	excepting both rates relative equally balanced.

$$\text{Averaged Error} = \frac{1}{n} \sum_{i=1}^n |P_i - P_{\hat{i}}|$$

The average error is the sum of all errors divided by the number of samples.

Date of birth class - 5

$$\text{Averaged Recall} = \frac{1}{n} \sum_{i=1}^n \frac{P_i}{P_i + F_{pi}}$$

The average recall is the average of all recalls.

2021

$$\text{Averaged F-measure} = \frac{1}{n} \sum_{i=1}^n \frac{P_i + R_i}{2P_i + R_i}$$

The average f-measure is the average of all f-measures.

F-measure

Write a short note on the following

a) Data Normalization

Normalization is used to scale the data of an attribute so that it falls in a smaller range, such as -1.0 to 1.0 or 0.0 to 1.0. It is generally useful for classification algorithms.

Methods of data normalization

### D) Decimal scaling method

In normalization by moving the decimal point of values of the data to normalize the data by this technique we divide each value of the data by the maximum absolute value of data the data value  $v_i$  of data. Then by using the formula below

$$v'_i = v_i / 10^j$$

Where  $j$  is the smallest integer such that  $v'_i$  is in the range [0, 1].

Ex - At the input data is : - 10 201 , 301 , 401

$10^{\pm 109'105}$

To normalize the above data,

*i* Maximum absolute values given data

cm) 87.01

Sol 2 I divided the given data by 1000 i.e.,  
The ungrouped data is:

0.01, 0.201, 0.301, -0.401, 0.501, 0.601, 0.701

2) min - max Normalization for this, technology data normalization is major needed

To this technology some prefer modern

affine linear transformation is

- the original data members is .  
- the original data is fetched & each value is .  
value from data is fetched & each value is .  
value from data is fetched & each value is .

Duplicated from - MARK (AN) - New - M10 (A) +

$$N^i = V - \text{Min}^i(A) \quad \text{Need}_i = \text{Min}^i(A)$$

where  $A_{ij}$  is the following  
 $\min(A)$  and  $\max(A)$  are the minimum & maximum  
 min (A) \* max (A) are the respectively  
 absolute value of a respectively entry in data  
 absolute value of each entry is called  
 max. The max value of each entry is called  
 max. The max value of each entry is called  
 min (A) is the min value of each entry  
 max (A) is the max value of each entry  
 min - max (A), max - min (A) respectively  
 max - min of the range respectively

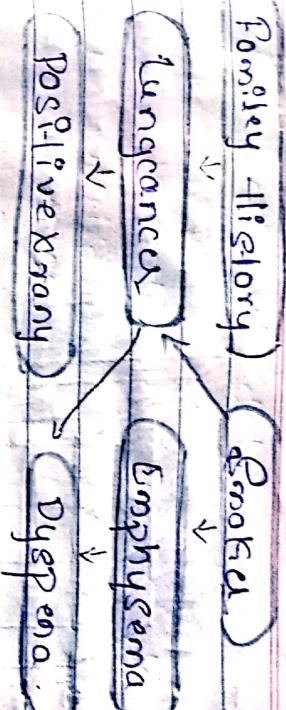
min NDBS = 1000 mg sulphur

5) 2 - score Normalization values and unit the  
To this a standard deviation of the

Bayesian belief networks

21

- \* Bayesian Belief Networks specify joint conditional probabilities. They are also known as Belief Revision Function.
- \* In belief networks all nodes cause conditional independencies to be defined between subsets of variables.
- \* It provides a graphical model of causal relations on which learning can be performed.
- \* We can use a revised Bayes' network for classification.
- \* Two components that define a Bayesian Belief Network



Q) A set of conditional probability table :-

The conditional probability table for the values of the variable lung cancer (L.C.) showing each possible combination of the values of its parent nodes. Family history (F.H.) known (G) is as follows

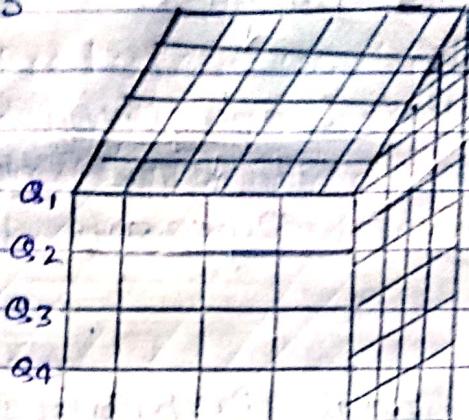
	F.H.S	F.H.-S.	F.H.S	F.H.-S
L.C.	0.8	0.5	0.7	0.1
L.C.	0.2	0.5	0.3	0.9

C) OLAP

ONLINE ANALYTICAL PROCESSING  
OLAP is a category of software that allows user to analyze information from different points of view

OLAP databases are divided into one or more cubes. The cubes are designed in such a way that creating & viewing reports become easy.

## OLAP cube :-



PG books shoe clothes

At the core of the OLAP concept is an OLAP cube. The OLAP cube is a data structure optimized for very quick data analysis.

The OLAP cube consists of numeric facts called measures which are categorized by dimensions. OLAP cube also called a hypercube. OLAP contains multidimensional data with data usually obtained from a different unrelated source.