

## I - ASSIGNMENT

(Start Writing From Here)

1. Define Data and Information. Explain Kinds of Data with an example.

A) Data :-

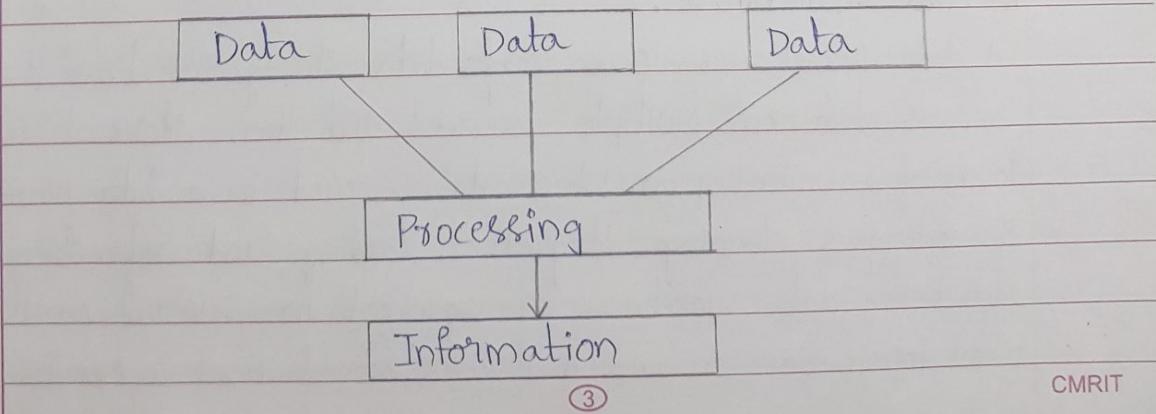
Data is a raw and unorganized fact that required to be processed to make it Meaningful. Data can be simple at the same time unorganized unless it is organized. Data comprises facts, observations, perceptions, numbers, characters, symbols, image etc..

Data is always interpreted, by a human or machine, to derive Meaning. Data contains numbers, statements and characters in a raw form.

Information :-

Information is a set of data which is processed in a meaningful way according to the given requirement. Information is processed, structured, or presented in a given context to make it Meaningful and useful. It improves the reliability of the data.

It helps to ensure undesirability and reduces uncertainty.



Types of data that can be Mined are :-

1. Data stored in the database :-

A database is also called a database Management System or DBMS. Every DBMS stores data that are related to each other in a way or the other. It also has a set of software programs that are used to Manage data and provide easy access to it and Managing different types of data access, such as shared, distributed and concurrent.

A relational db has tables that have different names, attributes and can store rows or records of large data sets. Every record stored in a table has a unique Key. Entity-relationship model is created to provide a representation of a relational db that features entities and the relationships that exist between them.

Ex :-

customer (cust\_ID, name, address, age, occupation, annual\_income, credit\_information, category, ...)

item (item-ID, brand, category, type, price, Supplier, cost...)

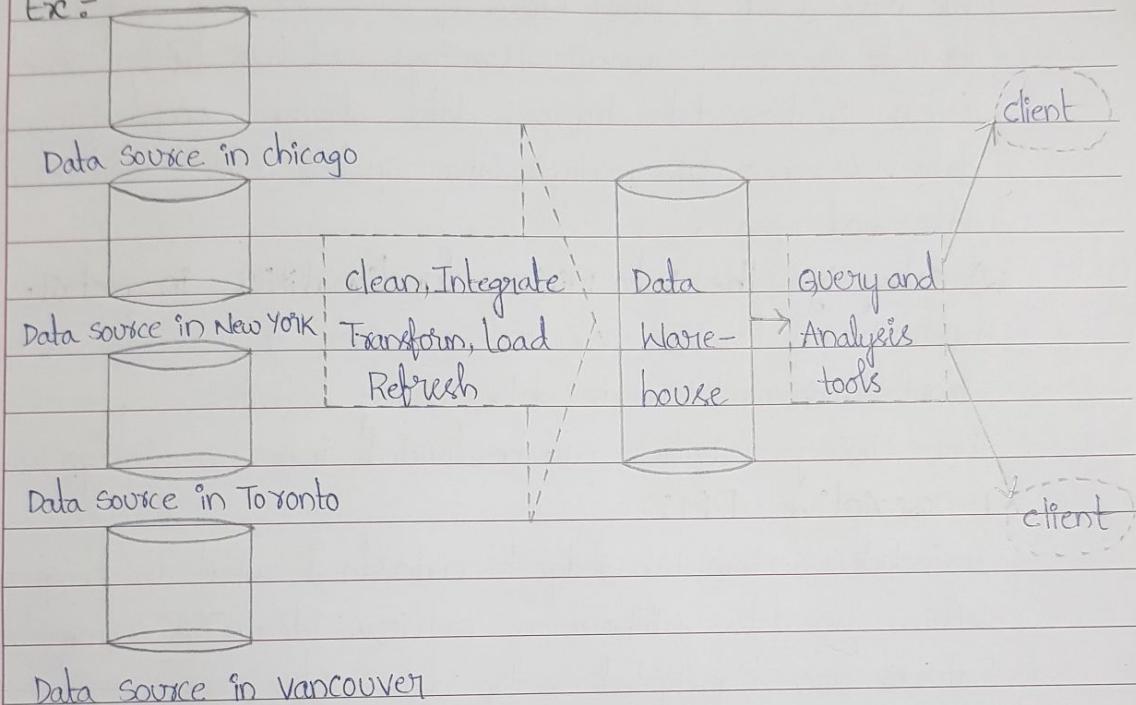
branch (branch-ID, name, address, ...)

2. Data Warehouse :-

A data Warehouse is a single data storage location that collects data from Multiple sources and then stores it in the form of a unified plan. When data is stored in a data Warehouse, it undergoes cleaning, integration, loading and refreshing. Data stored in a data Warehouse is organized in several parts. If you want information on data that was stored 6 or 12 Months back.

you will get it in the form of a summary.

Ex :-



Data source in Vancouver

### 3. Transactional data :

Transactional database stores record that are captured as transactions. These transactions include flight booking, customer purchase, click on a website, and others. Every transaction record has a unique ID. It also lists all those items that made it a transaction.

Ex :-

	trans-ID	list-of-item-IDs
	T100	I1, I3, I8, I6
	T200	I2, I8
	...	...

#### 4. other types of data :

We have a lot of other types of data as well that are known for their structure, semantic meanings, and versatility. Here are a few of those data types : data streams, engineering design data, sequence data, graph data, spatial data, multimedia data, and more.

#### 2) categorize the Data Mining functionalities in detail.

A) Data Mining functions are used to define the trends or correlations contained in data mining activities.

Data Mining activities can be divided into 2 categories :

##### 1. Descriptive DM:

It includes certain knowledge to understand what is happening within the data without a previous idea. The common data features are highlighted in the data set.

Ex: count, average etc..

##### 2. Predictive DM:

It helps developers to provide unlabeled definitions of attributes. Based on previous tests, the SW estimates the characteristics that are absent.

Ex: Judging from findings of a patient's Medical examinations.

##### DM functionality :

###### 1. class / concept Descriptions :

classes or definitions can be correlated with results. In simplified descriptive and yet accurate ways, it can be helpful to define individual groups and concepts. These class or concept definitions

are referred to as class/concept descriptions.

⇒ Data characterization :

This refers to the summary of general characteristics or features of the class that is under the study.

⇒ Data Discrimination :

It compares common features of class which is under study.

The output of this process can be represented in many forms.

Ex: bar charts, curves, pie charts.

2. Mining Frequent Patterns, Associations and correlations :

Frequent patterns are nothing but things that are found to be most common in the data.

There are different kinds of frequency that can be observed in the dataset.

⇒ Frequent item set :

This applies to a number of items that can be seen together regularly for ex: Milk and sugar.

⇒ Frequent Subsequence :

This refers to the pattern series that often occurs regularly such as purchasing a phone followed by a back cover.

⇒ Frequent Substructure :

It refers to the different kinds of data structures such as trees and graphs that may be combined with the itemset or subsequence.

Association analysis :

The process involves uncovering the relationship b/w data

and deciding the rules of the association. It is a way of discovering the relationship b/w various items. Ex: It can be used to determine the sales of items that are frequently purchased together.

#### correlation Analysis :

correlation is a Mathematical technique that can show whether and how strongly the pairs of attributes are related to each other. Ex: Heighted people tend to have more weight.

#### 3) classification and prediction :

##### ⇒ classification

It is the procedure of discovering a model that represents and distinguishes data classes or concepts, for the objective of being able to use the model to predict the class of objects whose class label is anonymous. The derived model is established on the analysis of a set of training data.

##### ⇒ Prediction

It defines predict some unavailable data values or pending trends. An object can be anticipated based on the attribute values of the object and attribute values of the classes. It can be a prediction of missing numerical values or increase /decrease trends in time-related information.

#### 4) clustering :

It is similar to classification but the classes are not predefined. The classes are represented by data attributes. It is unsupervised learning. The objects are clustered or grouped, depends on the

principle of Maximizing the intraclass similarity and Minimizing the interclass similarity.

### 5) Outlier analysis :-

Outliers are data elements that cannot be grouped in a given class or cluster. These are the data objects which have multiple behaviour from the general behaviour of other data objects. The analysis of this type of data can be essential to mine the knowledge.

### 6) Evolution analysis :-

It defines the trends for objects whose behaviour changes over some time.

## 3) Elaborate the different data preprocessing techniques

Data preprocessing is a data mining technique which is used to transform the raw data in a useful and efficient format.

Techniques involved are

### 1) Data cleaning :-

The data can have many irrelevant and missing parts. To handle this part, data cleaning is done. It involves handling of missing data and noisy data.

a) Missing data : This situation arises when some data is missing in the data. It can be handled in various ways.

i) Ignore the tuples : This approach is suitable only when the dataset we have is quite large and multiple values are missing within a tuple.

ii) Fill the Missing values : You can choose to fill the missing

values manually, by attribute mean or the most probable value.

b) Noisy data : It is a meaningless data that can't be interpreted by machines. It can be generated due to faulty data collection handled by.

1) Binning Method : It works on sorted data in order to smooth it. The whole data is divided into segments of equal size and then various methods are performed to complete the task.

2) Regression : Data can be made smooth by fitting it to a regression function. Regression used may be linear or multiple.

3) Clustering : It groups the similar data in a cluster.

2) Data Transformation :

It transforms data in appropriate forms suitable for Mining process.

1) Normalization : It is done in order to scale the data values in a specified range (-1.0 to 1.0 or 0.0 to 1.0)

2) Attribute Selection : In this, new attributes are constructed from the given set of attributes to help the Mining process.

3) Discretization : It is done to replace the raw values of numeric attribute by interval levels or conceptual levels.

4) Concept Hierarchy Generation : Attributes are converted from lower level to higher level in hierarchy.

3) Data reduction : It aims to increase the storage efficiency and reduce data storage and analysis costs.

1) Data cube Aggregation : Aggregation operation is applied to data for the construction of the data cube.

2) Attribute subset selection : The highly relevant attributes

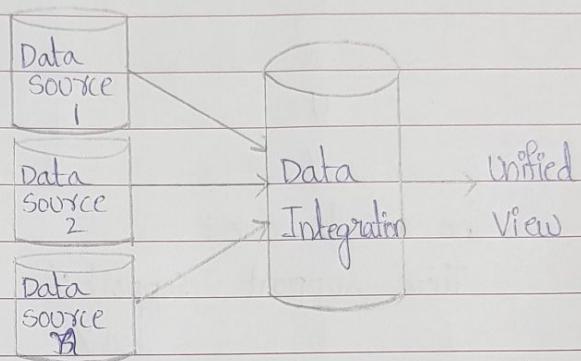
should be used, rest all can be discarded.

3) Numerosity Reduction: This enable to store the model of data instead of whole data.

4) Dimensionality reduction: This reduce the size of data by encoding Mechanisms.

4) Data Integration:

It combines data from multiple sources into a coherent data source. These sources may include multiple databases.



ii) Apply the following rules on a database has 5 transactions.

let min sup = 60% and min items bought conf = 80%.

TID	Items bought
T100	{M,O,N,K,E,Y}
T200	{D,O,N,K,E,Y}
T300	{M,A,K,E}

TID	Items bought
T400	{M,U,C,K,Y}
T500	{C,O,O,K,I,E}

a) Find all frequent item sets using Apriori

b) list all of the strong association rules (with sup & conf)

A) i. The database scanned once to generate the 1 itemset. The total ID is 5. Minimum support is 60%.

Item	Support	Support %
A	1	20%
C	2	40%

D	1	20%
E	5	100%
I	1	20%
K	4	80%
M	3	60%
N	2	40%
O	4	80%
U	1	20%
Y	3	60%

The results are:

Item support Support%  $\Rightarrow$  2. The database scanned once to generate the 2 Itemsets.

E	5	100%
K	4	80%
M	3	60%
O	4	80%
Y	3	60%

Item	Support	Support%	Item	Support	Support%
E,K	4	80%	M,Y	2	40%
E,M	3	60%	O,Y	3	60%
E,O	4	80%	The result is:		
E,Y	3	60%	E,K	4	80%
K,M	2	40%	E,M	3	60%
K,O	3	60%	E,O	4	80%
K,Y	2	40%	E,Y	3	60%
M,O	2	40%	K,O	3	60%

3. The database scanned once again to generate 3 itemsets. Result is:

Item	Support	Support%
E,K,O	3	60%
E,O,Y	3	60%

The itemsets (E,K,O) and (E,O,Y). So, there are 6 possible Association rules.

Association rules of {E,K,O}:

$$R_1: E \rightarrow K,O = 60\% < 80\%$$

$$R_2: K \rightarrow E,O = 75\% < 80\%$$

$$R_3: O \rightarrow E,K = 75\% < 80\%$$

$$R_4: E,K \rightarrow O = 75\% < 80\%$$

$$R_5: K,O \rightarrow E = 100\% > 80\%$$

$$R_6: E,O \rightarrow K = 75\% < 80\%$$

Association rules of {E,O,Y}:

$$R_7: E \rightarrow O,Y = 60\% < 80\%$$

$$R_8: K \rightarrow E,O = 100\% > 80\%$$

$$R_9: O \rightarrow E,Y = 75\% < 80\%$$

$$R_{10}: E,O \rightarrow Y = 75\% < 80\%$$

$$R_{11}: O,Y \rightarrow E = 100\% > 80\%$$

$$R_{12}: E,Y \rightarrow O = 100\% > 80\%$$

Rules R<sub>5</sub>, R<sub>8</sub>, R<sub>11</sub>, R<sub>12</sub> are strong association rules

∴ customers who purchase any of the 2 products from E,K,O would want to purchase the other one and all the customers who purchase 2 items from E,O,Y also want to purchase the other one.

4) What is FP-Growth tree? Explain FP-Growth tree algorithm with an Example.

A) Frequent pattern tree is a tree-like structure that is made with initial itemsets of the db. Each node of FP tree represents an item of the itemset. The rootnode represents null while the lower nodes represent the itemsets.

## FP Algorithm steps

The FP growth method lets us find the Frequent patterns without candidate generation.

- #1 First scan the database to find the occurrences of the itemsets in db. i.e support count.
- #2 secondly, construct the FP tree. create root node which is null.
- #3 scan the db again and examine the transactions. The itemset with the Max count is taken at top, the next set with lower count and so on. tree constructed in descending order.
- #4 After ordering if any itemset of this transaction is already present in another branch then this would share a common prefix to the root. common itemsets are linked.
- #5 count of IS is incremented as it occurs in transactions.
- #6 Mine the created FP tree. Traverse the path which is conditional pattern
- #7 construct a conditional FP tree, which is formed by a count of itemsets in the path. The itemsets Meeting threshold support are considered in the conditional FP tree.
- #8 Frequent patterns are generated from this.

Ex:- Support threshold = 50%, confidence = 60%.

Table 1

Transaction	List of items	$\text{Min-Supp} = 3 = 0.5 \times 6$
T <sub>1</sub>	i <sub>1</sub> , i <sub>2</sub> , i <sub>3</sub>	
T <sub>2</sub>	i <sub>2</sub> , i <sub>3</sub> , i <sub>4</sub>	
T <sub>3</sub>	i <sub>4</sub> , i <sub>5</sub>	
T <sub>4</sub>	i <sub>1</sub> , i <sub>2</sub> , i <sub>4</sub>	
T <sub>5</sub>	i <sub>1</sub> , i <sub>2</sub> , i <sub>3</sub> , i <sub>5</sub>	
T <sub>6</sub>	i <sub>1</sub> , i <sub>2</sub> , i <sub>3</sub> , i <sub>4</sub>	(14)

1. Count of each item and sorting them in descending order

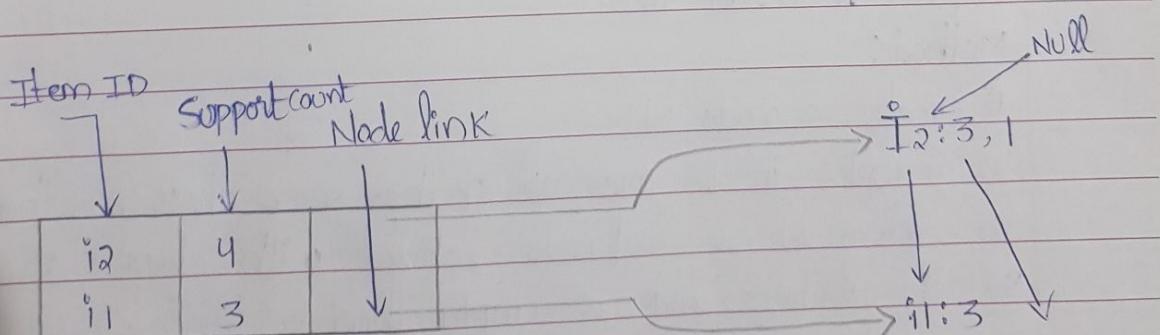
Item	count	Item	count
i1	4	i2	5
i2	5	i1	4
i3	4	i3	4
i4	4	i4	4
i5	2		

3. Build FP tree

Null  
 ↘ ↘  
 i2:5 i4:1  
 ↘ ↘ ↓  
 i1:4 i3:1 i5:1  
 ↓ ↓  
 i3:3 i4:1  
 ↓ ↘  
 i5:1 i4:1

4. Mining of FP-tree

Item	conditional Patterns base	conditional FP-tree	Frequent Patterns generated
i4	{i2, i1, i3:1}	{i2:2, i1, i4:2}	
	{i2, i3:1}	i3:2	{i3, i4:2}
			{i2, i3, i4:2}
i3	{i2, i1:3}	{i2:4, i1:3}	{i2, i3:4}
	{i2:1}		{i1, i3:3}
			{i2, i1, i3:3}
i1	{i2:4}	{i2:4}	{i2, i1:4}



5. Explain Decision tree induction algorithm for classification.

Discuss the usage of information gain.

A machine researcher named J. Ross Quinlan in 1980 developed a decision tree Algorithm known as ID3 (Iterative dichotomiser). later he presented C4.5, which was the successor of ID3.

In this algorithm, there is no backtracking; the trees are constructed in a top-down recursive divide and conquer manner.

Decision tree induction algorithm

Algorithm: Generate\_decision\_tree.

Input:

- ⇒ Data partition D, which is a set of training tuples and their associated class labels.
- ⇒ attribute\_list, the set of candidate attributes
- ⇒ Attribute\_selection\_method, a procedure to determine the splitting criterion that "best" partitions the data tuples into individual classes. This criterion consists of a splitting\_attribute and possibly, either a split-point or splitting\_subset.

Output: A decision tree.

Method:

- 1) Create a node N;
- 2) if tuples in D are all of the same class, c, then
- 3, return N as leaf node labeled with class c;
- 4) if attribute\_list is empty then
- 5) return N as a leaf node labeled with majority class in D;
- 6) apply Attribute\_selection\_Method(D, attribute\_list) to find the

"best" splitting criterion;

- 7) label node  $N$  with splitting criterion;
- 8) if splitting attribute is discrete valued and Multiway splits allowed then // not restricted to binary trees.
- 9) attribute list  $\leftarrow$  attribute list - splitting attribute,
- 10) for each outcome  $j$  of splitting criterion
- 11) let  $D_j$  be the set of data tuples in  $D$  satisfying outcome  $j$
- 12) if  $D_j$  is empty then
- 13) attach a leaf labeled with the majority class in  $D$  to node  $N$ ;
- 14) else attach the node returned by Generate decision tree  $(D_j, \text{attribute\_list})$  to node  $N$ ;
- 15) return  $N$ ;

Usage of information gain:

Information gain helps to determine the order of attributes in the nodes of a decision tree. we can use  $IG_i$  to determine how good the splitting of nodes in a decision tree.

# I - ASSIGNMENT

(Start Writing From Here)

- 1) What is Data mining? Explain Data mining preprocessing Techniques.

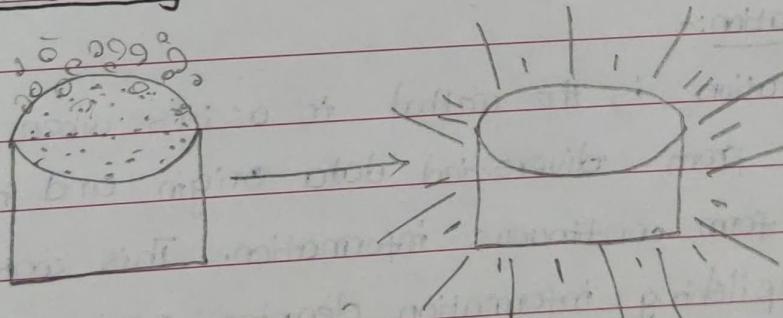
Data Mining:- is the process of sorting through large data sets to identify patterns and relationships that can help solve business problems through data analysis.

Preprocessing:-

- 1) Incompleteness (lacking attributes)
- 2) Noise (Deviate from actual result)
- 3) Inconsistency (Discrepancies in code).

Forms of pre-processing:-

- A) Data cleaning:-



The data can have many irrelevant and missing parts. To handle this part, data cleaning is done. It involves handling of missing data, noisy data etc.

- (a). missing data:-

This situation arises when some data is missing in the data. It can be handled in various ways.

Some of them are:

### 1. Ignore the tuples:-

This approach is suitable when some data is missing in the data. It can be handled in various ways. dataset we have is quite large and multiple values are missing within a tuple.

### 2) Fill the missing values:-

There are various ways to do this task. You can choose to fill the missing values manually, by attribute mean or the most probable value.

## (B) Noisy Data :-

### Data Integration:-

Data integration is the method to assists when the information is collected from diversified data origin and information is merging to form continuous information. This continuous information after accomplishing information cleaning assists in analysis and data preparation.

### Data Transformation:-

Data Transformation is assisted to modify the raw information into a definite design according to the requirement of the set.

### Normalization:-

In the Normalization method, numerical information is modified into the identified range.

### Aggregation:-

Aggregation can be specified from the conversation itself, this process assists to merge the characteristics into one.

### Generalization:-

In generalization, a lower level feature is modifying to a higher state.

### Data Reduction:-

After the conversion and scaling of information duplication that is insignificance within the information is taking away and effectively organizes the information during data preparation.

Q) List out Data mining task primitives?

A) A data mining task can be specified in the form of a data mining query, which is input to the data mining system. A data mining query is defined in terms of data mining task primitives.

\* These primitives allow the user to interactively communicate with the data mining system during discovery to direct the mining process.

1. The set of task-relevant data to be mined:-

This specifies the portions of the database (or) the set of data in which the user is interested. This includes the database

attributes over data warehouse dimensions of interest.

- \* In a relational database, the set of task-relevant data can be collected via a relational query involving operations like selection, projection, join and aggregation.
- \* The data collection process results in a new data relational called the "initial data relation".

② The kind of knowledge to be mined:-

This specifies the data mining functions to be performed, such as characterization, discrimination, association (or) correlation analysis, classification, prediction, clustering, outlier analysis, (or) evolution analysis.

3- The background knowledge to be used in the discovery process.

This knowledge about the domain to be mined is useful for guiding the knowledge discovery process and evaluating the patterns found. Concept hierarchies are a popular form of background knowledge, which allows data to be mined at multiple levels of abstraction.

- ° Rolling Up - Generalization of data:- Allow to view data at more meaningful and explicit abstractions and makes it easier to understand. It compresses the data, and it would require fewer input/output operations.

o Drilling Down - Specialization of data:- concept values replaced by lower-level concepts. Based on different user viewpoints, there may be more than one concept hierarchy for a given attribute or dimension.

u) The interestingness measures and thresholds for pattern evaluation:-

Different kinds of knowledge may have different interesting measures. They may be used to guide the mining process or, after discovery, to evaluate the discovered patterns. For example, interesting measures for association rules include support and confidence. Rules whose support and confidence values are below user-specified thresholds are considered uninteresting.

→ The to make pattern evaluation.

- o simplicity
- o certainty (confidence)
- o utility
- o Novelty

5) The expected representation for visualizing the discovered patterns:-

This refers to the form in which discovered patterns are to be displayed, which may include rules, tables, cross tabs, charts, graphs, decision trees, cubes, or other visual representations.

\* users must be able to specify the forms of presentation to be used for displaying the discovered patterns. Some representations may be better suited than others for particular kinds of knowledge.

3) Illustrate with an Apriori Algorithm for the given dataset below.

TID	list of items
001	milk, dal, sugar, bread
002	Dal, sugar, wheat, jam
003	milk, bread, curd, paneer
004	wheat, paneer, dal, sugar
005	milk, paneer, bread
006	wheat, dal, paneer, bread

A) 1. Generating 1<sup>st</sup> frequent itemset pattern

Items	Support count
Milk	3
dal	4
sugar	3
paneer	4
bread	4
wheat	3
Jam	1 X
curd	1 X

< 1

1 or < 2 will be neglected.

Items	Support count
Milk	3
dal	4
Sugar	3
paneer	4
bread	4
wheat	3

&lt;1

2) Generating 2<sup>nd</sup> frequent item set pattern.

Items	Support count
milk, dal	1 <input checked="" type="checkbox"/>
milk, sugar	1 <input checked="" type="checkbox"/>
milk, paneer	2
milk, bread	3
milk, wheat	0 <input checked="" type="checkbox"/>
dal, sugar	3
dal, paneer	2
dal, bread	2
dal, wheat	3
sugar, paneer	1 <input checked="" type="checkbox"/>
sugar, bread	1 <input checked="" type="checkbox"/>
sugar, wheat	2
paneer, bread	3
paneer, wheat	2
bread, wheat	1 <input checked="" type="checkbox"/>

1 or &lt;2 will be neglected.

Items	Support count
milk, paneer	2
milk, bread	3
Dal, Sugar	3
Dal, paneer	2
Dal, bread	2
Dal, wheat	3
Sugar, wheat	2
paneer, bread	3
paneer, wheat	2

3) Generating 3rd frequent item set pattern

Items	Support count
Milk, paneer, bread	2
Milk, paneer, dal, sugar	0
Milk, paneer, dal	0
Milk, dal, paneer, bread	0
Milk, paneer, dal, wheat	0
Milk, paneer, sugar, wheat	0
Milk, paneer, bread	2
Milk, paneer, wheat	0
Milk, bread, dal, sugar	1
Milk, bread, dal, paneer	0
Milk, bread, dal	1
Milk, bread, dal, wheat	0

milk, bread, sugar, wheat	0	X
milk, bread, paneer	2	
milk, bread, wheat, paneer	0	X
Dal, sugar, paneer	1	X
Dal, sugar, bread	1	X
Dal, sugar, wheat	2	
Dal, sugar, wheat	2	
Dal, sugar, paneer, bread	0	X
Dal, sugar, paneer, wheat	1	X
Dal, paneer, bread	1	X
Dal, paneer, wheat	2	
Dal, bread, wheat	1	X
Dal, bread, sugar, wheat	0	X
Dal, bread, paneer	1	X
Dal, bread, paneer, wheat	1	X
Dal, wheat, sugar	2	
Dal, wheat, paneer, bread	1	X
Dal, wheat, paneer, wheat	2	
Sugar, wheat, paneer, bread	0	X
Sugar, wheat, paneer	1	X
paneer, bread, wheat	1	X

1 or 2 will be neglected.

Items	Support count
milk, paneer, bread	3
Dal, sugar, wheat	3
Dal, paneer, wheat	3

Dal, paneer, wheat is the most associative frequent item set.

- v) Apply the following rules on the database has five transactions  
 Let min sup = 60%. And min items bought conf = 80%.

T <sub>100</sub>	Items bought
T <sub>100</sub>	{M, O, N, K, E, Y}
T <sub>200</sub>	{D, O, N, K, E, Y}
T <sub>300</sub>	{M, A, K, E}
T <sub>400</sub>	{M, O, C, K, Y}
T <sub>500</sub>	{C, D, O, K, I, E}

- (i) Apply FP-Growth Algorithm to the following transactional data to find frequent items.  
 (ii) List all frequent itemsets with their support count.

## FP-Growth Algorithm.

1) Generating 1st frequent itemset pattern:-

TID	Datasets.
T <sub>100</sub>	{M, O, N, K, E, Y}
T <sub>200</sub>	{O, N, K, G, Y}
T <sub>300</sub>	{N, A, K, E}
T <sub>400</sub>	{M, U, C, K, Y}
T <sub>500</sub>	{C, O, A, K, I, E}

2)	Item	support count.	support%.
	M	3	60%
	O	3	60% <del>max-</del>
	N	2	40%
	K	5	100%
	E	4	80%
	Y	3	60%
	O	1	20%
	A	1	20%
	U	1	20%
	C	2	40%
	I	1	20%

8) The results are :-

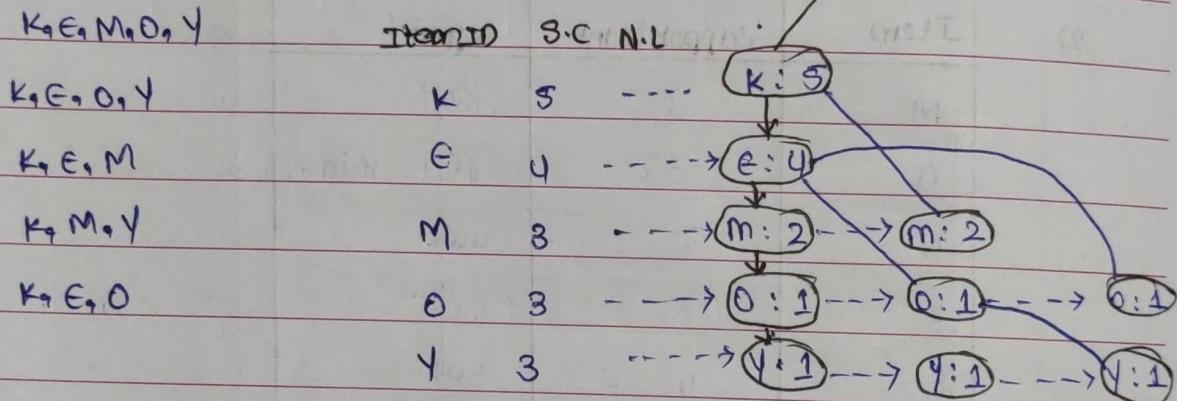
Item	support	support%.
E	5	100%
K	4	80%
M	3	60%
O	4	80%
Y	3	60%

⇒ The database scanned once to generate the 2 Item sets



frequent pattern.

u) Ordered Item set:- s) F-P tree:-



6) Mining of FP-tree:-

Item	conditional pattern base	conditional F.P tree.	Generated F.P.
Y	[Σ KEMO: 13, Σ KE: 13, Σ KM: 13]	[Σ K: 33]	<K, Y : 3>
O	[Σ KEM: 13, Σ KE: 23]	[Σ KE: 33]	<K, O : 3> <E, O: 3> <O, E, K: 3>
M	[Σ KE: 23, Σ K: 13]	[Σ K: 33]	<M, K: 3>
E	[Σ K: 13]	[Σ K: 13]	<E, K: 3>
K	-	-	-

s)

Explain naive Bayes algorithm with an example.

A)

Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.

- \* The Naive Bayes algorithm is comprised of 2 words Naive and Bayes.

### Baye's Theorem:-

- \* Baye's Theorem is a simple mathematical formula used for calculating conditional.
- \* Conditional probability is a measure of the probability of an event occurring given that another event has (by assumption, presumption, assertion, or evidence) occurred.

probability of B occurring given  
evidence A has already occurred

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

probability of A occurring

↑

P(B)

probability of A occurring  
given evidence B has  
already occurred

probability of B occurring

Example:-

person	Height	Weight	Foot size.
male	6.00	180	12
male	5.92	190	11
male	5.58	170	12
male	5.92	165	10
Female	5.00	100	6
Female	5.50	150	8
Female	5.42	130	7
Female	5.75	150	9.

$$P(\text{male}) = 4/8 = 0.5$$

$$P(\text{female}) = 4/8 = 0.5$$

Male:-

$$\text{Mean (Height)} = \frac{(6+5.92+5.58+5.92)}{4} = 5.855$$

$$\text{Variance (Height)} = \frac{\sum (x_i - \bar{x})^2}{n-1} = 0.035055$$

Sex	(height) mean	(height) variance	(weight) mean	(weight) variance	(footsize) mean	(footsize) variance
Male	5.855	0.035055	176.25	122.92	11.25	0.91667
Female	5.4175	0.097225	132.5	0558.33	7.5	1.6667

$$P(\text{male}) = 4/8 = 0.5$$

$$P(\text{female}) = 4/8 = 0.5$$

New instance to be classified is:-

Sex	Height	weight	foot size.
sample	6	130	8.

$$P(H|M) = \frac{1}{\sqrt{2 \cdot 3.142 \cdot 0.035033}} * e^{-\frac{(6 - 5.855)^2}{2 \cdot 0.035033}} = 1.5789.$$

$P(W M) = 5.9881e^{-6}$	$P(H F) = 2.2346 e^{-1}$	$P(FS F) = 2.8669 e^{-1}$
$P(FS M) = 1.3112 e^{-3}$	$P(W F) = 1.6789 e^2$	

posterior (male) =  $\frac{p(cm) * p(H|m) * p(w|m) * p(FS|m)}{\text{Evidence.}}$

$$\Rightarrow 0.5 * 1.5789 * 5.9881 e^{-6} * 1.3112 e^{-3} = 6.198 e^{-9}$$

posterior (female) =  $\frac{p(F) * p(N|F) * p(W|F) * p(FS|F)}{\text{Evidence.}}$

$$\Rightarrow 0.5 * 2.2346 e^{-1} * 1.6789 e^2 * 2.8669 e^{-1}$$

$$\Rightarrow 5.377 e^{-4}$$

# I - ASSIGNMENT

(Start Writing From Here)

- 1) what is KDD? Explain KDD steps with a neat diagram.

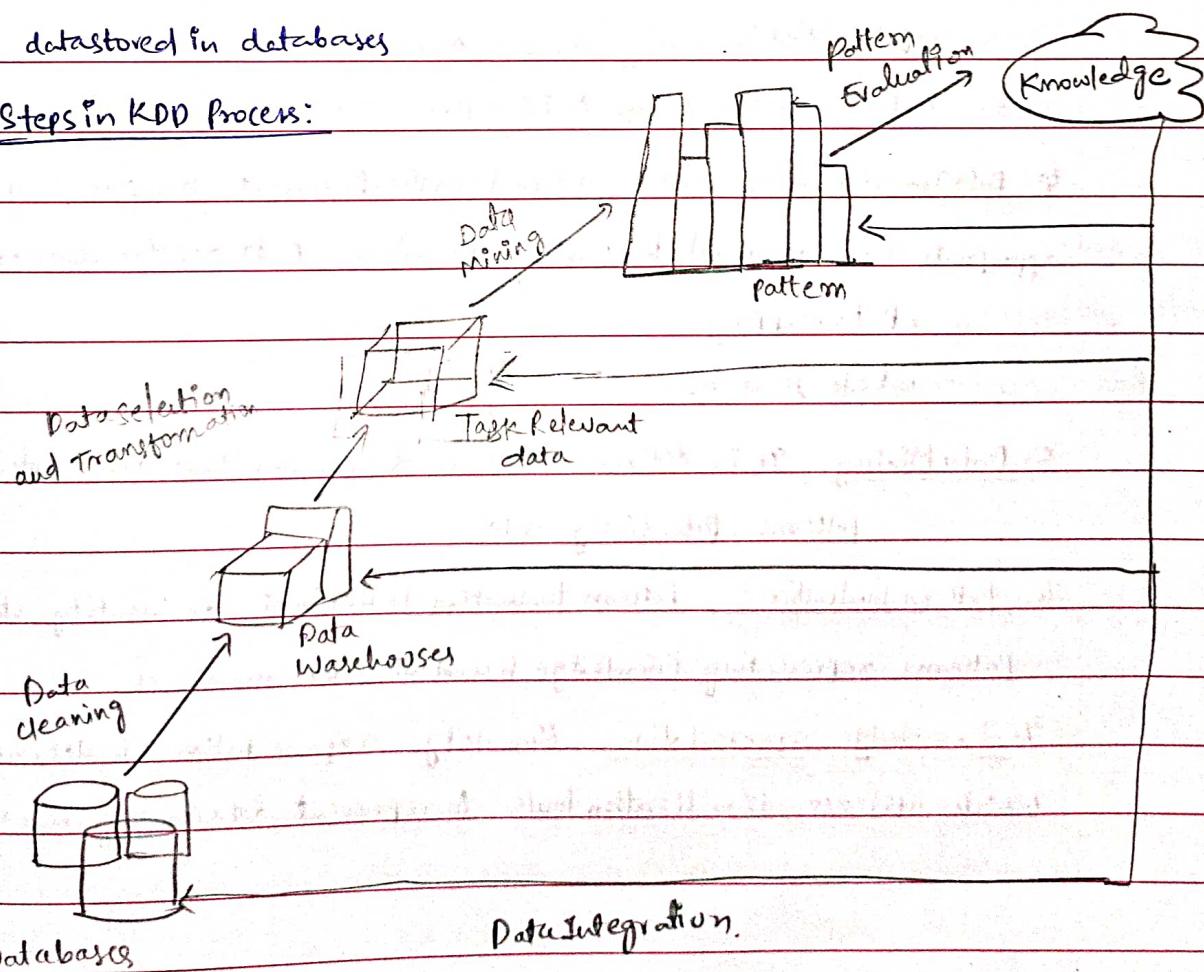
Data Mining - Knowledge Discovery in Databases (KDD).

why we need Data Mining?

volume of information is increasing everyday that we can handle from business transactions, Scientific data, Sensor data, Pictures, Videos etc. So, We need a system that will be capable of extracting essence of information available and that can automatically generate report.

\* Data Mining also known as knowledge Discovery in Databases, refers to the nontrivial extraction of implicit, Previously Unknown and Potentially useful Info. from data stored in databases

Steps in KDD Process:



1. Data Cleaning: Data cleaning is defined as removal of noisy and irrelevant data from collection.

- cleaning in case of missing values
- cleaning noisy data, where noise is a random or variance error.

2. Data Integration: Data integration is defined as heterogeneous data from multiple sources combined in a common source

- Data Integration Using Data Migration tools.
- Data Integration Using Data Synchronization tools

3. Data Selection: Data selection is defined as the process where data relevant to the analysis is derived and retrieved from the data collection.

- Data Selection Using Neural Network
- Data Selection Using Decision Trees.

4. Data Transformation: It is defined as the process of transforming data into appropriate form required by mining procedure. Data transformation is a two step process

- Data Mapping
- Code generation

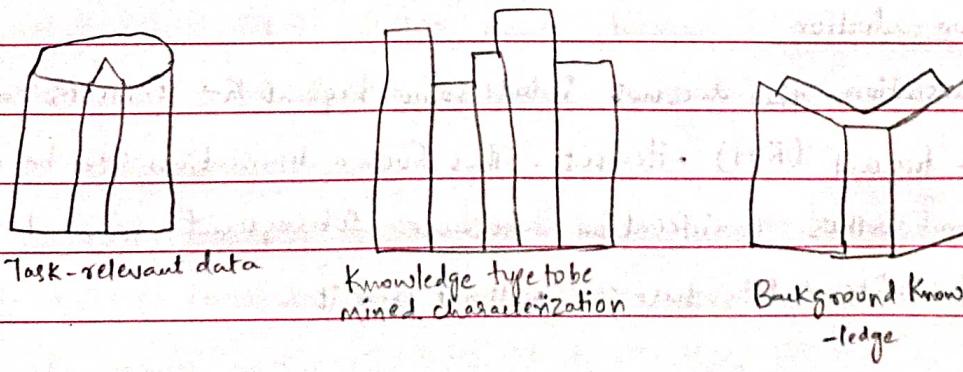
5. Data Mining: It is defined as clever techniques that are applied to extract patterns potentially useful.

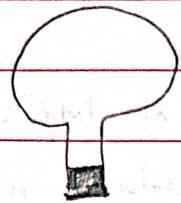
6. Pattern Evaluation: Pattern evaluation is defined as identify strictly increasing patterns representing knowledge based on given measures.

7. Knowledge representation: Knowledge representation is defined as technique which utilizes visualization tools to represent data mining results.

2) List out issues in Data Mining and Data Mining task Primitives.

- \* Data Mining Task can be Specified in the form of datamining Query which
- \* A data mining query is defined in terms of datamining task primitives.
- \* These primitives allow the user to interactively communicate with the data mining system during discovery to direct the Mining process or examine the findings from different angles or depths.
- \* The data Mining Primitives Specify the following
  1. Set of task-relevant data to be Mined.
  2. Kind of Knowledge to be Mined.
  3. Background Knowledge to be Used in the discovery process.
  4. Interestiness measures and thresholds for Pattern Evaluation
  5. Representation for Visualizing the discovered Patterns.
- \* A data mining Query language can be designed to incorporate these primitives, allowing users to interact with data Mining Systems flexibly. Having a data mining Query Language provides a foundation on which user-friendly graphical interfaces can be built.





Pattern Interestingness  
measures Simplicity.

gfg
dfdf
dfdfdf

Visualization of discovered Patterns, Rules,  
tables, reports, charts, graphs.

- 3) Explain improving of the efficiency of Apriori Algorithm.

These are some variations of the Apriori algorithm that have been projected that target developing the efficiency of the original algorithm which are as follows

#### 1) The hash-based technique :

A hash-based technique can be used to decrease the size of the Candidate K-itemsets,  $C_K$ , for  $K > 1$ . For instance, when scanning each transaction in the database to create the frequent 1-itemsets,  $L_1$  from the Candidate 1-itemsets in  $L_1$ , it can make some 2-itemsets for each transaction, hash them into several buckets of a hash table structure, and increase the equivalent bucket counts.

#### 2) Transaction reduction :

A transaction that does not include some frequent  $K$ -itemsets cannot include some frequent  $(K+1)$ -itemsets. Thus such a transaction can be marked or deleted from further consideration because subsequent scans of the database for  $j$ -itemsets, where  $j > K$  will not need it.

3) Partitioning: - A partitioning technique can be used that needed 2 database scans to mine the frequent itemsets. It includes 2 phases involving in phase 1, the algorithm subdivides the transactions of D into n non-overlapping partitions. If the minimum support count for a partition is  $\text{min\_sup} \times$  the number of transactions in that partition.

4) Sampling: The fundamental idea of the Sampling is to select a random sample of the given data D, and then search for frequent itemsets in S rather than D. In this method, it can trade off some degree of accuracy against efficiency. The sample size of S is such that the search for frequent itemsets in S can be completed in Main Memory, and therefore only one scan of the transactions in S is needed overall.

5) Explain naive Bayes algorithm with an example.

Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.

\* The Naive Bayes algorithm is comprised of 2 words Naive and Baye's, which can be described as:

\* Naive: It is called Naive because it assumes that the occurrence of a certain feature is independent of the occurrence of other features.

\* Baye's: It is called Baye's because it depends on the principle of Baye's Theorem.

Baye's Theorem:

\* Baye's Theorem is also known as Baye's Rule or Baye's law, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability.

\* The formula for Bayes theorem is given as

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

where,

- \*  $P(A|B)$  is Posterior Probability
- \*  $P(B|A)$  is Likelihood Probability.
- \*  $P(A)$  is Prior Probability
- \*  $P(B)$  is Marginal Probability.

Working of Naive Bayes classifier:

1. Convert the given dataset into frequency table.
2. Generate likelihood table by finding the probabilities of given features.
3. Now, use Bayes theorem to calculate the Posterior Probability.

problem: If the weather is Sunny, then Player Should Play or not?

Solution: To Solve this, first Consider the below dataset.

	OUTLOOK	Play
0	Rainy	Yes
1	Sunny	Yes
2	overcast	Yes
3	overcast	Yes
4	Sunny	No
5	Rainy	Yes
6	Sunny	Yes
7	overcast	Yes

8	Rainy	No
9	Sunny	No
10	Sunny	Yes
11	Rainy	No
12	Overcast	Yes
13	Overcast	Yes

Frequency table for the weather conditions

Weather

Weather	Yes	No
Overcast	5	0
Rainy	2	2
Sunny	3	2
Total	10	5

Likelihood table Weather Condition:

Weather	No	Yes	
Overcast	0	5	$5/14 = 0.35$
Rainy	2	2	$4/14 = 0.29$
Sunny	2	3	$5/14 = 0.35$
All	$4/14 = 0.29$	$10/14 = 0.71$	

Applying Baye's Theorem:

$$P(\text{Yes}|\text{sunny}) = P(\text{sunny}|\text{Yes}) \cdot P(\text{Yes}) \\ P(\text{sunny}).$$

$$P(\text{Sunny} | \text{Yes}) = 3/10 = 0.3$$

$$P(\text{Sunny}) = 0.35$$

$$P(\text{Yes}) = 0.7$$

$$\text{So, } P(\text{Yes} | \text{Sunny}) = 0.3 \times 0.7 / 0.35 = 0.60$$

$$P(\text{No} | \text{Sunny}) = \frac{P(\text{Sunny} | \text{No}) \times P(\text{No})}{P(\text{Sunny})}$$

$$P(\text{Sunny} | \text{No}) = 2/4 = 0.5$$

$$P(\text{No}) = 0.29$$

$$P(\text{Sunny}) = 0.35$$

$$\text{So, } P(\text{No} | \text{Sunny}) = \frac{0.5 \times 0.29}{0.35} = 0.41$$

So we can see from above calculation that  $P(\text{Yes} | \text{Sunny}) > P(\text{No} | \text{Sunny})$ .

Ans FP-Growth Algorithm is Frequent pattern growth algorithm to improve apriori method. FP growth algorithm represents the database in the form of a tree called as FP Tree.

FP Tree: It is a tree like data structure that is made with initial itemsets of the database.

The main purpose is to mine the most frequent pattern.

### 1. Count of each item

Item	Count
$i_1$	2
$i_2$	2
$i_3$	5
$i_4$	5
$i_5$	5
$i_6$	5
$i_7$	4
$i_8$	2

### 2. Sort the itemset in descending order

Item	min-support count
$i_3$	5
$i_4$	5
$i_5$	5
$i_6$	5
$i_7$	4

$i_8$	2	
$i_2$	2	
$i_1$	2	

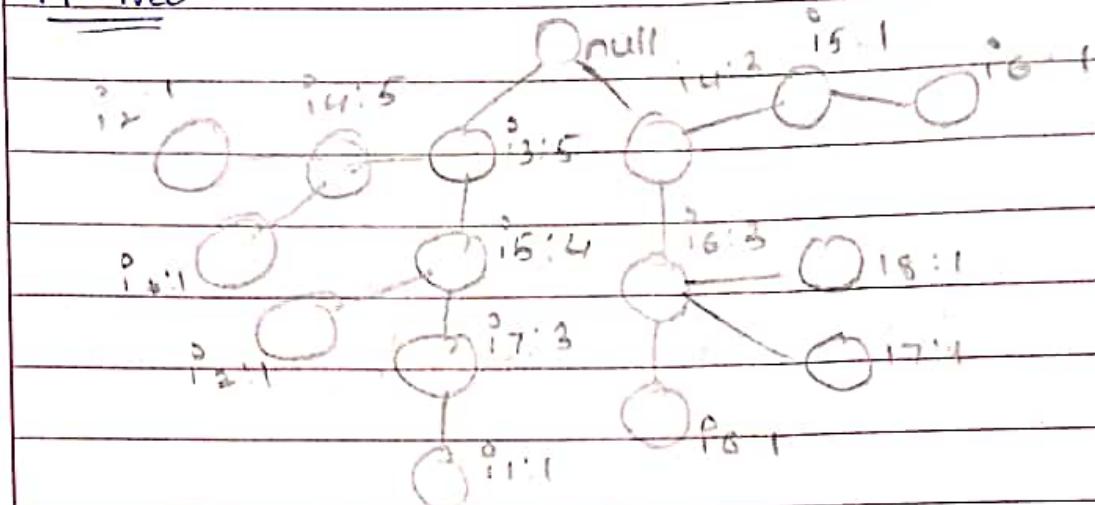
FP-tree A frequent pattern set ( $L$ ) is built which contains all the elements whose frequency is greater than or equal to minimum support.

$$L = \{i_3: 5, i_4: 5, i_5: 5, i_6: 4, i_7: 4\}$$

Transaction ID	Items	Ordered - Item Set
I <sub>1</sub>	{i <sub>3</sub> , i <sub>5</sub> , i <sub>7</sub> }	{i <sub>3</sub> , i <sub>5</sub> , i <sub>7</sub> } ✓
I <sub>2</sub>	{i <sub>4</sub> , i <sub>6</sub> , i <sub>8</sub> }	{i <sub>4</sub> , i <sub>6</sub> , i <sub>8</sub> } ✓
I <sub>1</sub>	{i <sub>3</sub> , i <sub>5</sub> , i <sub>7</sub> }	{i <sub>3</sub> , i <sub>5</sub> , i <sub>7</sub> } ✓
I <sub>9</sub>	{i <sub>7</sub> , i <sub>5</sub> , i <sub>1</sub> }	{i <sub>5</sub> , i <sub>7</sub> } ✓
I <sub>2</sub>	{i <sub>4</sub> , i <sub>6</sub> , i <sub>7</sub> }	{i <sub>4</sub> , i <sub>6</sub> , i <sub>7</sub> } ✓
I <sub>1</sub>	{i <sub>2</sub> , i <sub>3</sub> , i <sub>4</sub> }	{i <sub>3</sub> , i <sub>4</sub> } ✓
I <sub>3</sub>	{i <sub>4</sub> , i <sub>5</sub> , i <sub>6</sub> }	{i <sub>4</sub> , i <sub>5</sub> , i <sub>6</sub> } ✓
I <sub>7</sub>	{i <sub>8</sub> , i <sub>6</sub> , i <sub>1</sub> }	{i <sub>6</sub> } ✓
I <sub>8</sub>	{i <sub>5</sub> , i <sub>3</sub> , i <sub>2</sub> }	{i <sub>3</sub> , i <sub>5</sub> } ✓
I <sub>9</sub>	{i <sub>3</sub> , i <sub>4</sub> , i <sub>6</sub> }	{i <sub>3</sub> , i <sub>4</sub> , i <sub>6</sub> } ✓

Inserting each element or item into tree

### FP-tree



5Q Explain naïve Bayes algorithm with an example.

Ans Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.

\* The Naïve Bayes algorithm is comprised of 2 words Naïve and Bayes.

### Baye's Theorem

It is used to determine the probability of a hypothesis with prior knowledge.

It depends on conditional probability

# I - ASSIGNMENT

(Start Writing From Here)

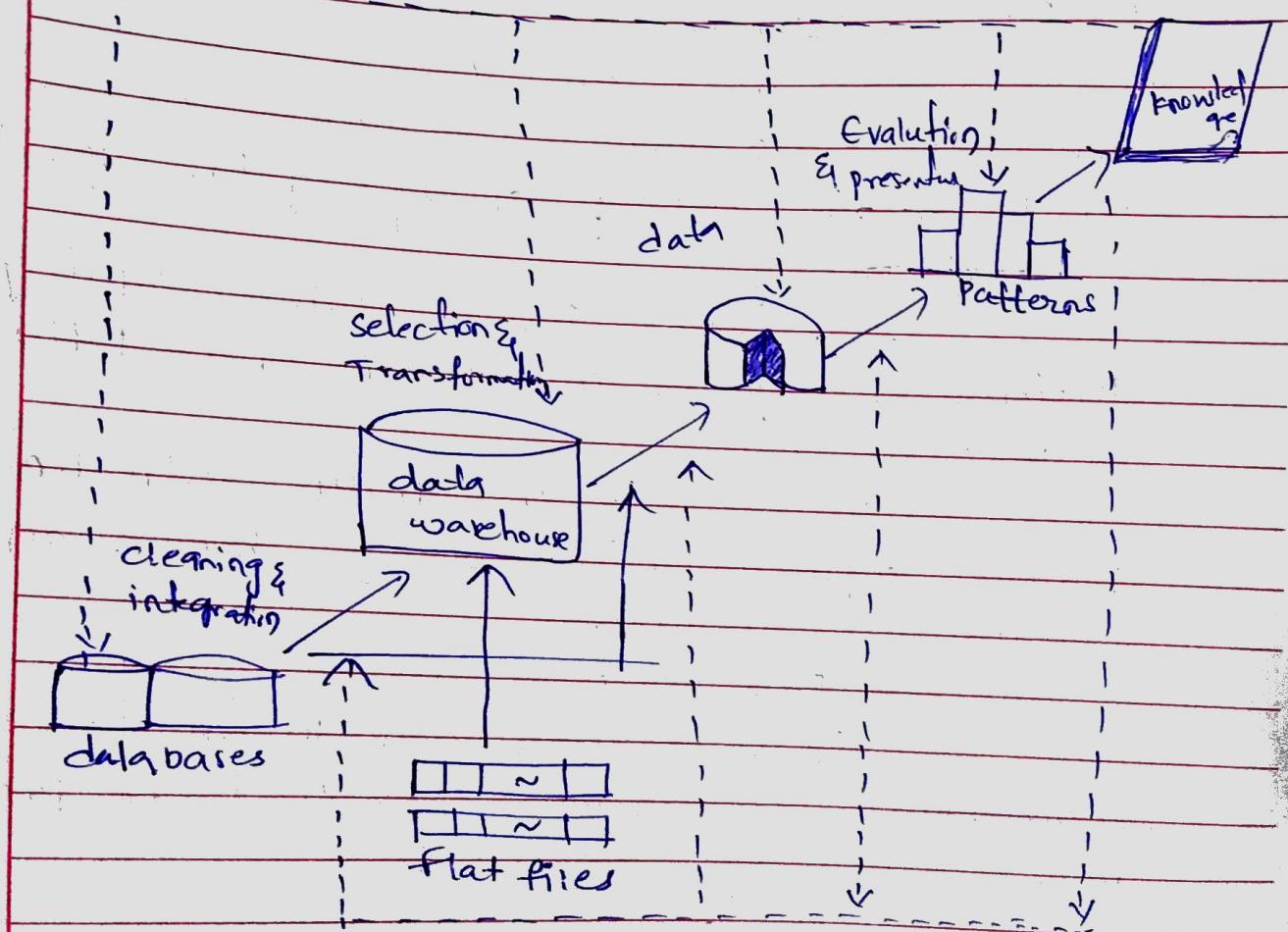
①

Ans

What is KDD? Explain KDD steps with a neat diagram.

Data Mining is a fast growing field also known as knowledge discovery from data, or KDD or knowledge mining from data, knowledge extraction data/patterns analysis, data archaeology and data dredging.

The knowledge discovery process is an iterative sequence of the following steps:



## 1. Data cleaning

To remove noise and inconsistent data

## 2. Data integration

Where multiple data sources may be combined.

## 3. Data selection

Where data relevant to the analysis task are retrieved from the data base

## 4. Data Transformation

Where data are transformed and consolidated into forms appropriate for mining by performing summary or aggregation operations

## 5. Data mining

An essential process where intelligent methods are applied to extract data patterns

## 6. Pattern evolution

To identify the truly interesting patterns representing knowledge based on interestingness measures

## 7. Knowledge presentation

Where visualization and knowledge representation techniques are used to present mined knowledge to users.

(2)

Write a short notes on

- a) Data Cleaning
- b) Data Transformation

An

### Q) Data Cleaning

it is a crucial process in data mining.  
it carries an important part in the building  
of a model.

Data cleaning can be regarded as the process  
needed, but very often neglects it.

Data Quality is the main issue in quality  
information Management.

Data quality problems occur anywhere in  
information systems.

These problems are solved by data cleaning.

Data Cleaning is fixing or removing  
incorrect, corrupted, incorrectly formatted  
duplicate, or incomplete data with a  
dataset.

### Steps of Data cleaning

- 1) Remove duplicate
- 2) fix structural errors
- 3) filter unwanted outliers
- 4) handle Missing data

## 5) validate and QA

### Methods of data cleaning

- 1) ignore the tuples
- 2) fill the missing value
- 3) Binning
- 4) Regression
- 5) Clustering

### Characteristics of data cleaning

- Accuracy
- coherence
- Validity
- uniformity
- Data verification
- clean data back flow

## b) Data Transformation

it is the process of converting data from one format to another, typically from that format of a source system into the required format of a destination system.

Data Transformation is a component of most data integration and data Management tasks, such as data wrangling and data warehousing.

The data are transformed in ways that are ideal for mining the data.

The data transformation involves steps that are:

1. Smoothing
2. Aggregation
3. Discretization
4. Attribute construction
5. Generalization
6. Normalization
  - i) Min-Max Normalization
  - ii) Z-score Normalization,
  - iii) Decimal scaling .

③

Illustrate with an A-priori Algorithm for the given dataset below:

TID	list of items
001	Milk, dal, sugar, bread
002	Dal, sugar, wheat, jam
003	Milk, bread, curd, paneer
004	wheat, paneer, dal, sugar
005	Milk, paneer, bread

006 wheat, dal, Paneer, bread.

Ans: 1. Generating 1<sup>st</sup> frequent itemset pattern

items	support count
Milk	3
dal	4
Sugar	3
Paneer	4
bread	4
wheat	3
Jam	D X
Curd	D X

C1

1 or < 2 will be neglected.

items	support count
Milk	3
dal	4
Sugar	3
Paneer	4
bread	4
wheat	3

L1

## 2) Generating 2<sup>nd</sup> frequent item set pattern.

items	Support count
Milk, da	1 X
Milk, sugar	1 X
Milk, paneer	2
Milk, bread	3
Milk, wheat	0 X
dal, sugar	3
dal, paneer	2
dal, bread	2
dal, wheat	3
Sugar, paneer	1 X
Sugar, bread	1 X
Sugar, wheat	2
paneer, bread	3
paneer, wheat	2
bread, wheat	1 X C1

1 or < 2 will be neglected

items	support count
Milk, Paneer	2
Milk, bread	3
dal, sugar	3
dal, Panner	2
dal, bread	2
dal, wheat	3
sugar, wheat	2
Panna, bread	3
Panner, wheat	2

L1

3) Generating 3rd frequent item set patterns.

items	support count
Milk, Paneer, bread	2
Milk, Paneer, dal, sugar	0) X
Milk, Panner, dal	0) X
Milk, dal, Panner, bread	0) X
Milk, panner, dal, wheat	0) X
Milk, Panna, sugar, wheat	0) X
Milk, Panner, bread	2
Milk, Panner, wheat	0) X
Milk, bread, dal, sugar	1) X
Milk, bread, dal, panner	0) X

Milk, bread, dal	1) X
Milk, bread, dal, wheat	0) X
Milk, bread, sugar, wheat	0) X
Milk, bread, panner	2
Milk, bread, wheat, panner	0) X
dal, sugar, panner	5) X
dal, sugar, bread	1) X
dal, sugar, wheat	2.
dal, sugar, wheat	2
dal, sugar, panner, bread	0) X
dal, sugar, panner, wheat	1) X
dal, panner, bread	1) X
dal, panner, wheat	2
dal, panner, sugar, wheat	1) X
dal, panner, bread	1) X
dal, panner, wheat	2
dal, bread, wheat	5) X
dal, bread, sugar, wheat	0) X
dal, bread, panner	1) X
dal, bread, panner, wheat	1) X
dal, wheat, sugar	2.
dal, wheat, panner, bread	1) X
dal, wheat, panner, wheat	2
Sugar, wheat, panner, bread	0) X
Sugar, wheat, panner	1) X

Panner bread wheat

c1

1 or <2 will be neglected

items	support count
Milk, panner, bread	3
dal, sugar, wheat	3
dal, panner, wheat	3.

dal, panner, wheat is the most associative frequent item set

- ④ What is correlation analysis? Explain Correlation Analysis from Association Analysis?
- Correlation analysis is a statistical method used to measure the strength of the linear relationship between two variables and compute their association.
- Correlation analysis calculates the level of change in one variable due to the change in the other.

A high correlation points to a strong relationship between the two variables, while a low correlation means that the variables are

weakly related.

Types of correlation analysis in Data mining.

1) Pearson correlation

$$\rho_{xy} = \frac{n \sum xy_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

2) Kendall rank correlation

$$\tau = \frac{nc - nd}{\frac{1}{2} n(n-1)}$$

$nc$  = no. of concordant

$nd$  = No. of discordant

3) Spearman rank correlation.

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

$\rho$  = Spearman rank correlation

$d_i$  = the diff b/w the ranks of corresponding variable,

$n$  = no. of observations

(5)

A

Explain Decision tree induction algorithm for classification. Discuss the usage of information gain.

Decision Tree is a Supervised learning method used in data mining for classification and regression methods. It is a tree that helps us in decision-making purposes. The decision tree creates classification or regression as a tree structure. It separates a data set into smaller subsets, and at the same time, the decision tree is steadily developed.

- The final tree is a tree with the decision nodes and leaf nodes
- A decision node has at least two branches
- The leaf nodes show a classification or decision.

### Key factors

#### Entropy:-

It refers to a common way to measure impurity. In the decision tree, it measures the randomness or impurity in data sets.

high entropy



less random

dataset



less entropy

information Gain:

Information Gain refers to the decline in entropy after the dataset is split. It is also called Entropy Reduction.

Building a decision tree is all about discovering attributes that return the highest data gain

The decision tree useful for:

- It enables us to analyze the possible consequences of a decision throughly
- it provides us a framework to measure the values of outcomes and the probability of accomplishing them
- it helps us to make the best decisions based on existing data and best speculations

Advantages of using decision-trees:

- A decision tree does not need scaling of information
- Missing values in data also do not influence the process of building a choice tree to any considerable extent.
- A decision tree model is automatic and simple to explain to the technical team as well as stakeholders.
- Compared to other algorithms, decision trees need less exertion for data preparation during pre-processing.
- A decision tree does not require a standardization of data.