

UNIT 4: Introduction to Data Mining

SQL → RDBMS ← simple query
Data Mining → Data Warehouse ← complex query

- Data Mining: access to data
- Data warehouse: managing the data
- Mining data from data warehouse is an advance version of extracting SQL data from RDBMS.

Problems with today's database management

The amount of data in computer files and database, business areas, banks is growing at phenomenal rate at the same time the users of these data are expecting much sophisticated information from them.

i.) The problem is how to store these data, how to manage these data and how to access these data in order to solve this problem, & the solution is data mining, and data warehouse.

When a marketing manager wants to predict what are the most important trends in customer behaviour.

what is the difference between SQL & data mining?

→

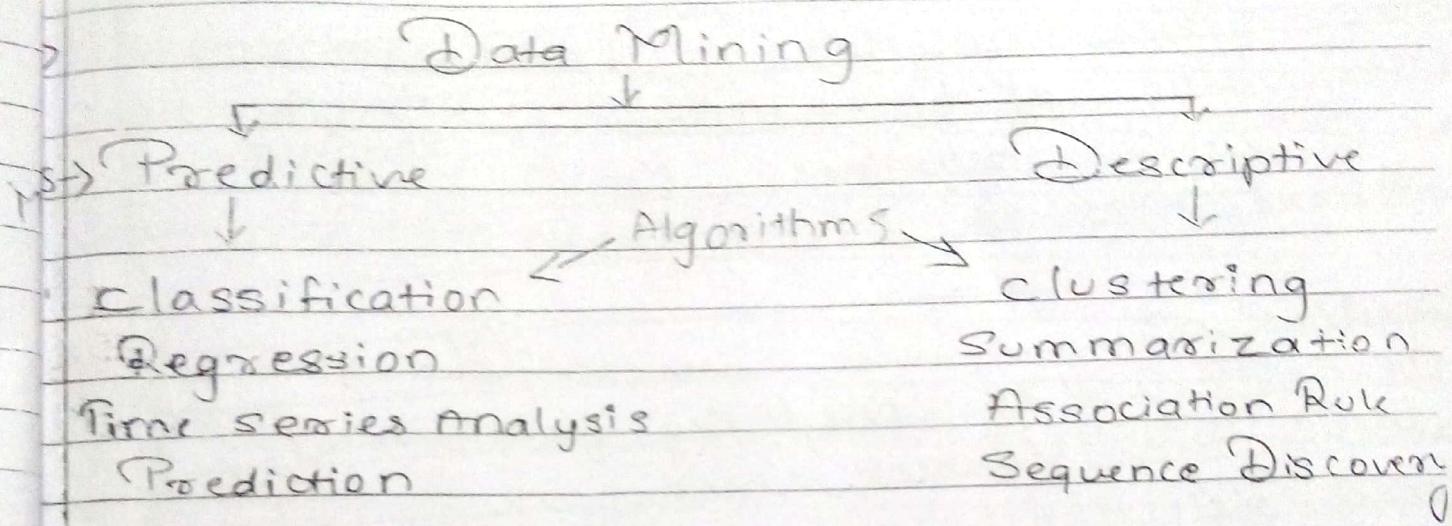
	SQL	Data Mining
Query	• SQL query are simple & queries are complex or and well formed might not be well-formed or precisely stated.	
Data	• The data accessed from the original operational dbs., which are not in consistent format. (RDBMS)	• Data accessed is different version from original operational db. The data have been cleansed and modified to better support the data mining.
Output	The output of SQL is subset of database.	• In data mining, you may vaguely know what you are looking for.

Definition of Data Mining:

Data mining is defined as finding hidden information in a database. It is also called as exploratory data analysis, data-discovery, and deductive learning.

Data mining is used to describe discovering a mining language knowledge from large amt. of

Data mining involves many different algorithm to accomplish task. It has two models: Predictive & Descriptive.]



Predictive Model:-

- It makes a prediction about values of data using known results from different data.
- Algorithms:
- Classification:-
- It maps data into predefined groups or classes.
Eg: an airport security screening is used to determine if passengers are terrorist or criminals.
- Regression:-
- It is used to map a data item to a real valued prediction variable.]
If an employee wants to reach a certain level of savings before his/her retirement

Saving will be based on current values and past values.

Time Series Analysis:

The value of an attribute is examined as it varies over time.

Share market / stock

Prediction:

It is used to predicting future state rather than current state.

flood detection.

Descriptive Method - Model:-

This model identifies patterns or relationships between data.

clustering:

It is called as segmentation.

It is similar as classification except that the groups are not predefined

The groups contain similar information.

Car Price	???
	???
	???
	???
	???

Age of customer

It abstracts representative info about database.

g: compare students of Mumbai and Pune university. It will use to estimate the intellectual level of student.

ii) Association Rule:

- It describes relationships among data.

Bread \leftrightarrow Butter

pencil \leftrightarrow eraser.

iii) Sequence Discovery:

- It is used to determine sequential pattern in data.

$\langle A, B, C, D, E \rangle$ } the pattern describes 70%.

$\langle A, B, C, E \rangle$ } users of page A, follow

$\langle A, B, C \rangle$ } page C. So we will add

link directly from Page A \rightarrow C.

16 KDD : Knowledge Discovery Process

- Data mining and KDD is often use inter-changeably.
- Data mining is the use of algorithm to extract the information and patterns derived by the KDD process.

F Defination of KDD:-

- It is the process of finding useful information and patterns in data.

Input to this process is data and O/P is useful information derived by the user.

The KDD process consists of 5 steps:-

i) Selection:-

The data needed for data mining process obtained from many different and heterogeneous data sources.

ii) Preprocessing:-

Erroneous data may be corrected or removed.

Ex:- Bank Branch
at Thane

Bank Branch
at kurla

CID	Name	Add	CID	Name	Address
1	ABC	Thane	1	ABC	kurla
2	XYZ	Bhandup	2	XYZ	Bhandup

customer 1 has got different address in different databases, which can either removed or corrected.

iii) Missing data must be supplied or predicted or removed.

Bank Branch at Thane

Bank branch at kurla

CID	Name	Address
1	ABC	Thane
2	XYZ	Bhandup

CID	Name	Address
1	ABC	-
2	XYZ	Bhandup

It shud be supplied
predicted
or removed

4) Remove false information :-

CID	Name	Add	Magazine Purchase Date	Magazine
1	ABC	Thane	1-1-1901	car
2	XYZ	Bhandup	1-1-1968	cosmetics

→ Here, on the date of 1901, the company is not existing that means false information is present in the database, which should be corrected or removed.

5) Remove Redundancy :-

→ The preprocessing is used to correct missing values, de-duplication, disambiguation before transferring data from operational System to data warehouse.

ii) Transformation :-

Data from different sources must be converted into a common format for processing.

- It is used to reduce the no. of possible data values.
- Either keep age or dob.
- Either keep no. of years of experience or date of joining.

iv) Data Mining :-

- It is used to apply many algorithms like

clustering, classification, association rule etc.
on the data present in data warehouse.

- Interpolation / Evaluation (Visualization Technique)
- How data mining results are presented to the user, is extremely important because usefulness of the result is dependent on it. For this visualization and GUI's strategy are used.

Visualization Technique include:

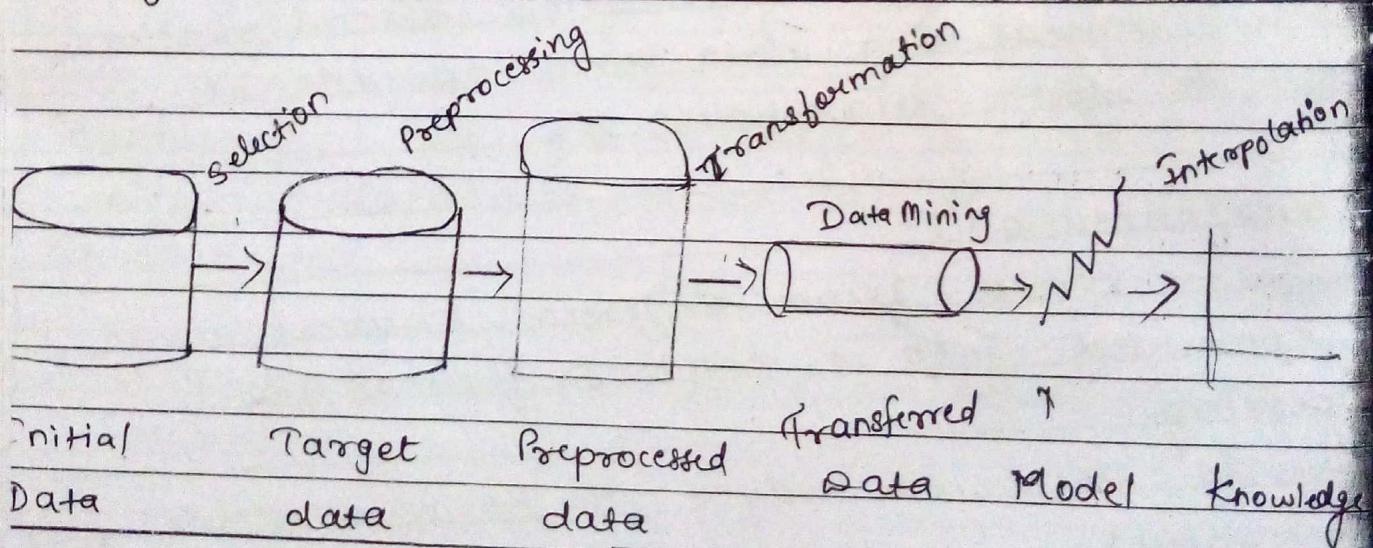
graphical

Icon based

pixel based

hierarchical

hybrid.



Data Mining Issues:-

Human Interaction

Visualization of results

Multimedia data

Missing data e.) Ease of use.

* Data Mining Data *

categorical

divide data

In specific
groups)

Nominal

Eg: marital
status,
gender

ordinal

height,
Credit
Score

Numerical

(represents

numerical
value)

interval

temperature

Ratio

mass,
energy.

- Data refers to collection of facts usually obtained as result of experiences and observations.

Data & Data mining refers as follows:-

* Categorical Data :- It represents labels of multiple classes used to divide a variable into specific groups.

i.) Nominal Data:- It contains obj. which are not measurements. Eg: marital status can be categorized as single / married / divorced.

ii.) Ordinal Data:- It contains code on object dat represent the rank order among them.

Eg: Credit score can be categorised as low/medium/high

* Numerical Data:- It represents numeric value of specific variables. ; Interval Data:-

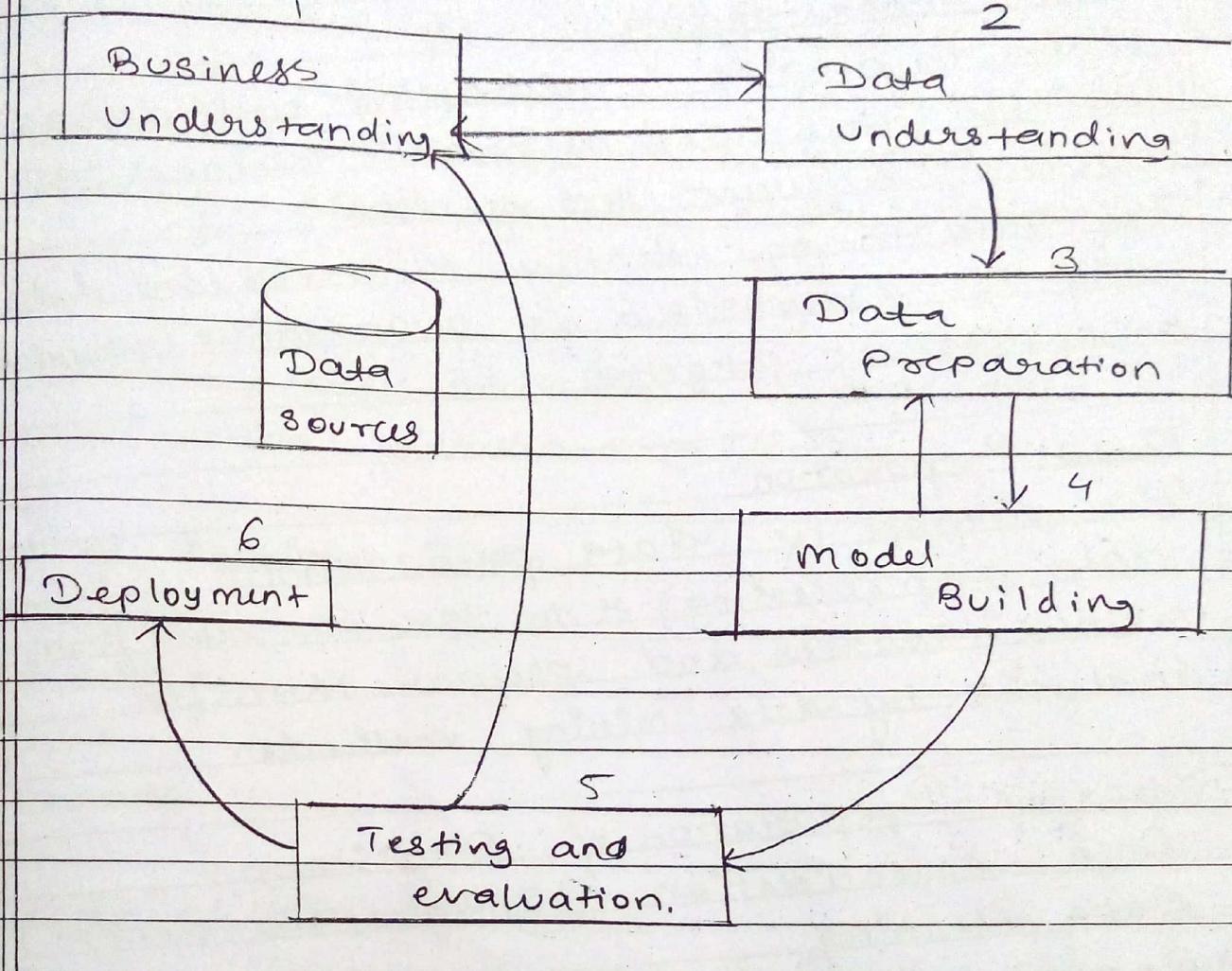
Eg: age, income, temperature. These are variables measured on interval scale. Eg. temp.

(ii) Ratio: It includes variables commonly found in the physical sciences and engineering.
Eg. mass, energy, time length.

Data Mining Applications:

- 1) Customer Relationship Mgmt. (CRM): The goal of CRM is to create one to one relationship with customers by developing & understanding their needs & wants.
- 2) Banking: It identifies base to maximize customer value by selling them products and services like atm machine, credit card, banking branches etc.
- 3) Insurance: It uses data mining technique to predict which customer are more likely to buy new policy.
- 4) Computer H/w & S/w: It uses data mining to detect & prevent comp. n/w security breach and identify potentially unsecure s/w products.
- 5) Entertainment Industry: Data Mining is used to analyse viewer's data to decide what programs to show during prime time & how to minimize returns by knowing where to insert advertisements.
- 6) Medicines: - Data Mining is used to discover the relationships between symptoms & illness. It predicts success rates of organ transplantation patients to develop better donor organ matching policies.

* Data Mining Process *



- Data Mining researchers have proposed several researches to minimize the chance of success in conducting data mining projects.

1) Business Understanding:-

- It is the thorough understanding of the managerial need for new knowledge; an explicit specification of the business objective regarding the study to be conducted.

- 2.) Data Understanding:-
- Data mining is addressing well defined business task and different business task require different sets of data
 - The data mining identify the relevant data which is categorized as quantitative (numerical) & qualitative (categorical)

3.) Data Preparation:

- The purpose of data preparation (data processing) is to take the data from various sources and prepare them for analysis by data mining methods.

4.) Steps of Preparation of Data:-

- 1.) Data consolidation (Integration / Selection)
- 2.) Data cleaning
- 3.) Data transformation
- 4.) Data Reduction

5.) Model Building:-

On prepared data, we will apply one of the data mining model to get the output.

6.) Testing and Evaluation:-

This step will test the output and evaluate it.

7.) Deployment:-

Once testing has been successful, we will deploy the output.

Association Rule

Date 1/1
Page 13
STUDY BUDDIES

The purchasing of one product, when another product is purchased, we present an association rule. It is used by retail stores, to assist in marketing advertising, floor placement and inventory control.

Association rules are used to show the relationships between data items.

* Uses of Association Rule:-

- a) Placement
- b) Advertising
- c) Sales
- d) Coupons

* Objective of Association Rule:

Increased sales & Reduced costs.

* Definition of Association Rule:-

Given a set of Items $I = \{I_1, I_2, \dots, I_m\}$ and a database of transaction $D = \{t_1, t_2, \dots, t_n\}$ where $t_i = \{I_{i1}, I_{i2}, \dots, I_{ik}\}$ and $I_{ij} \in I$, an association rule is an implication of the form $x \Rightarrow y$, where $x, y \subseteq I$ are sets of items called itemset and $x \cap y = \emptyset$.

eg:	Transactions	Items
		I_{00}
t_1		Bread, Jelly, Peanut Butter
t_2		I_{01} Bread, peanut Butter
t_3		Bread, milk, Peanut Butter
t_4		Beer, Bread
t_5		Beer, Milk

5/11/16

Monday

* Definition of Support (S)

Association Rule

- The support (S) for an A.R $x \Rightarrow y$ is the percentage of transactions in the database that contain

$$\left[\frac{x \cup y}{\text{Total no. of items}} \times 100 \right]$$

* Definition of confidence or strength (α)

- The confidence or strength for an association $x \Rightarrow y$ is the ratio of no. of transactions that contain $x \cup y$ to the no. of transactions that contain x

$$\left[\frac{x \cup y}{x} \times 100 \right]$$

Importance is measured by two features
 Support & Confidence

The selection of association rules is based on support and confidence.

Confidence measures the strength of the rule.
Support measures how often it should occur in the database.

$X \Rightarrow Y$	Support(s)	Confidence(α)
$X \quad Y$ Beer \Rightarrow Bread	$\frac{1}{5} \times \frac{20}{100} = 20\%$	$\frac{1}{2} \times \frac{50}{100} = 50\%$
$X \quad Y$ PeanutButter \Rightarrow Jelly	$\frac{1}{5} \times 100 = 20\%$	$\frac{1}{3} \times 100 = 33.3\%$
$X \quad Y$ Jelly \Rightarrow PeanutButter	$\frac{1}{5} \times \frac{20}{100} = 20\%$	$\frac{1}{1} \times 100 = 100\%$
$X \quad Y$ Jelly \Rightarrow Milk	$\frac{0}{5} \times 100 = 0\%$	$\frac{0}{1} \times 100 = 0\%$
$X \quad Y$ Bread \Rightarrow PeanutButter	$\frac{3}{5} \times \frac{20}{100} = 60\%$	$\frac{3}{4} \times \frac{25}{100} = 75\%$
$X \quad Y$ PeanutButter \Rightarrow Bread	$\frac{3}{5} \times 100 = 60\%$	$\frac{3}{3} \times 100 = 100\%$

* Lift:

- the definition of lift is defined as

$$\text{lift}(x \Rightarrow y) = \frac{\text{Supp}(x \cup y)}{\text{Supp}(x) \times \text{Supp}(y)}$$

or the ratio of the observed support to that of expected to that if x and y were independent.

* Conviction:

\rightarrow the conviction of a.r. is defined as

$$\text{conv}(x \Rightarrow y) = \frac{i - \text{supp}(x)}{1 - \text{conf}(x \Rightarrow y)}$$

: calculate lift of Bread \Rightarrow PeanutButter

$$= \frac{\text{Supp}(x \cup y)}{\text{Supp}(x) \times \text{Supp}(y)}$$

$$= \frac{\text{Supp}(\text{Bread} \Rightarrow \text{PeanutButter})}{\text{Supp}(\text{Bread}) \times \text{Supp}(\text{PeanutButter})}$$

$$= \frac{60\%}{80\% \times 60\%} = \frac{1}{80\%} = \frac{0.1}{0.01} = 10$$

$$= \frac{10}{80} = 1.2$$

ii) Conviction of Bread \Rightarrow Peanut Butter:

$$= 1 - \text{supp}(x)$$

$$1 - \text{conf}(x \Rightarrow y)$$

$$= 1 - 80\%$$

$$(= 75\%)$$

$$= \frac{1 - 0.8}{1 - 0.75} = \frac{0.2}{0.25} = 0.8 = 1.0$$

iii) Calculate lift & conviction for jelly & Peanut Butter

$$\text{lift} = \frac{\text{supp}(\text{jelly} \cup \text{PeanutButter})}{\text{supp}(\text{jelly}) + \text{supp}(\text{PeanutButter})}$$

$$= \frac{20\%}{20\% \times 60\%} = \frac{1}{60} = \frac{0.1}{6} = 1.6$$

$$\text{ii) Conviction: } \frac{1 - \text{supp}(\text{jelly})}{1 - \text{conf}(\text{jelly} \Rightarrow \text{PeanutButter})}$$

$$= \frac{1 - 20}{1 - 100} = \frac{1 - 0.2}{1 - 1} = \infty$$

* Frequent Item Set (Large Item Set)

→ A large frequent item set is an item set whose no. of occurrences is above a threshold value (s, α) we use the notation to indicate the complete set of large item set & l to indicate a specific large item set.

Eg: Suppose that the input support & confidence are $s = 30\%$ and confidence = 50% . respectively
 $L = \{\{\text{Beer}\}, \{\text{Bread}\}, \{\text{Milk}\}, \{\text{PeanutButter}\}, \{\text{Bread}\}, \{\text{PeanutButter}\}\}$

$$l = \{\text{Bread, PeanutButter}\}_{l \in L}$$

There are 2 non-empty subsets of l

i) $\{\text{Bread}\}$ & $\{\text{PeanutButter}\}$

Frequent itemset (large itemset with at least one item in l)

Support (Bread, PeanutButter)

Support (Bread)

$$= \frac{60}{80} - \frac{0.6}{0.8} = \frac{3}{4} = 0.75 = 75\%$$

ii) with 2nd one $l = \{\text{PeanutButter}\}$

Support (Bread, PB)

$$\frac{\text{Sup}(PB)}{80\%} = \frac{60\%}{60\%} = 1 = \frac{1 \times 100}{1} = 100\%$$

$$x \Rightarrow (l-x)$$

Bread \Rightarrow PB | PB \Rightarrow Bread.

- i) The A.R. rule is given by $x \Rightarrow (l-x)$ which is
 \Leftrightarrow Bread \Rightarrow Peanut Butter (x)
- confidence of $A.R. \{Bread \Rightarrow Peanut\ Butter\}$ is 75%,
 which is above threshold value, hence
 it is a valid association rule & is added
 to R \Leftarrow (Set of valid A.R.) = { $\{Bread \Rightarrow Peanut\ Butter\}$ }

- i) The A.R. is given by $x \Rightarrow (l-x)$ which is
 \Leftrightarrow Peanut Butter \Rightarrow Bread }
- confidence of $A.R. \{Peanut\ Butter \Rightarrow Bread\}$
 is 100%, which is above threshold value
 hence it is a valid A.R. & is added
 to R = { $\{Peanut\ Butter \Rightarrow Bread\}$ }

Apriori Algorithm:-

This algorithm is most well known A.R. & is used in most commercial product.

It uses following property any subset of large item set must be large.

Working of apriori algorithm :-

It is used to generate candidate item set of a particular size and then scan the database to count to see if they are large.

- During scan: (index 0...n) candidate of size i (C_i) are counted only those candidate that are large are use to generate candidates for the next pass i.e. L_i are use to generate C_{i+1}
- An item set is consider as a candidate only in all its subset are large to generate candidate of size $i+1$ joins are made of large item set found in the previous pass.

Eg : $S = \{30\} \cup \{50\}$.

Pass	Candidates	Large Itemset.
1	40 80 20 $\{\text{Beer}\}$, $\{\text{Bread}\}$, $\{\text{Jelly}\}$, $\{\text{milk}\}$, $\{\text{PB}\}$ 40 60	$\{\text{Beer}\}$ $\{\text{Bread}\}$ $\{\text{milk}\}$ $\{\text{PB}\}$
2	20 20 $\{\text{Beer, Bread}\}$ $\{\text{Beer, milk}\}$ $\{\text{Beer, PeanutButter}\}$ $\{\text{Bread, milk}\}$ $\{\text{Bread, PeanutButter}\}$ $\{\text{milk, PB}\}$	$\{\text{Bread, PB}\}$
3	$\{\text{Bread, PB}\}$ NULL	$\{\text{Bread, PB}\}$ \emptyset

Define Support and confidence to measure the strength of association and calculate Support and confidence for the association rule.

TID	Item
1	{Bread, Milk}
2	{Bread, Diaper, Beer, Eggs}
3	{Milk, Diaper, Beer, Cola}
4	{Bread, Milk, Diaper, Beer}
5	{Bread, Milk, Diaper, Cola}

$$i) \text{Support } (\{ \text{Milk, Diapers} \} \rightarrow \{ \text{Beer} \}) =$$

$$\frac{x \cup y}{\text{total no. of items}} \times 100$$

$$= \frac{2}{5} \times \frac{100}{20} = 2 \times 20 = 40\%$$

$$i) \text{Confidence } (\{ \text{Milk, Diapers} \} \rightarrow \{ \text{Beer} \}) =$$

$$\frac{x \cup y}{x} \times 100$$

$$= \frac{2}{3} \times \frac{100}{33.3}$$

$$= 2 \times 3.33 = 66.6\%$$

Q2) Eg: What do you mean by frequent itemset? Use Apriori algo to generate frequent itemset for the following by taking $S = 60\%$.

ass	candidates	Large Itemset
1.	$\begin{matrix} 80 & 80 & 60 \\ \{\text{Broad}\} \cup \{\text{Milk}\} \cup \{\text{Beer}\} & \{\text{Broad}\} \cup \{\text{Milk}\} & \{\text{Broad}\} \cup \{\text{Beer}\} \\ \{\text{Diapers}\} \cup \{\text{Eggs}\} \cup \{\text{Cold}\} & \{\text{Diapers}\} \cup \{\text{Eggs}\} & \{\text{Diapers}\} \cup \{\text{Cold}\} \\ 80 & 20 & 40 \end{matrix}$	$\{\text{Broad}\} \cup \{\text{Milk}\}$ $\{\text{Diaper}\} \cup \{\text{Beer}\}$
2.	$\begin{matrix} 60 & 60 \\ \{\text{Broad}, \text{Milk}\} \cup \{\text{Broad}, \text{Diaper}\} & \{\text{Broad}, \text{Milk}\} \cup \{\text{Diaper}\} \\ \{\text{Milk}, \text{Diaper}\} \cup \{\text{Broad}, \text{Beer}\} & \{\text{Milk}, \text{Diaper}\} \cup \{\text{Beer}\} \\ \{\text{Milk}, \text{Beer}\} \cup \{\text{Diaper}, \text{Beer}\} & \{\text{Diaper}, \text{Beer}\} \end{matrix}$	$\{\text{Broad}, \text{Milk}\} \cup \{\text{Broad}, \text{Diaper}\}$ $\{\text{Milk}, \text{Diaper}\} \cup \{\text{Diaper}, \text{Beer}\}$ $\{\text{Broad}, \text{Milk}, \text{Beer}\}$
3.	$\begin{matrix} 40 \\ \{\text{Broad}, \text{Milk}, \text{Diaper}\} \\ \{\text{Broad}, \text{Milk}, \text{Diaper}, \text{Beer}\} \\ \{\text{Broad}, \text{Diaper}, \text{Beer}\} \\ \{\text{Milk}, \text{Diaper}, \text{Beer}\} \end{matrix}$	\emptyset
4.	\emptyset	\emptyset

- * Issues with Apriori Algorithm:
 - Apriori algo seems simple only for small data sets.
 - In much larger data set especially those with huge amount of items present in low quantities and small amount of items present in big quantities.

The search and calculations becomes a computationally intensive process.

FP Tree [Frequent Pattern Tree].

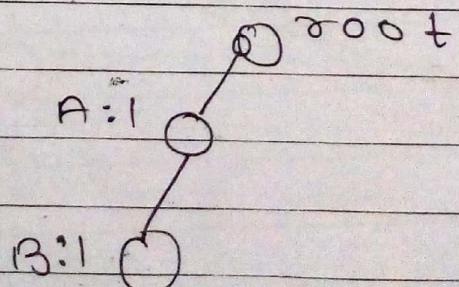
FP tree is a compact representation of the database once an FP tree has been constructed it uses a recursive divide & conquer approach to mine frequent item set.

Develop FP tree for following database :

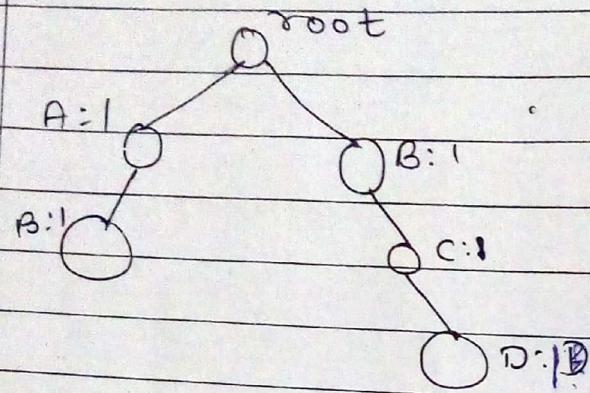
TID	Items
1	{A, B, C}
2	{B, C, D}
3	{A, C, D, E}
4	{A, D, E}
5	{A, B, C}
6	{A, B, C, D}
7	{B, C}
8	{A, B, C}
9	{A, B, D}
10	{B, C, E}

i.) After reading

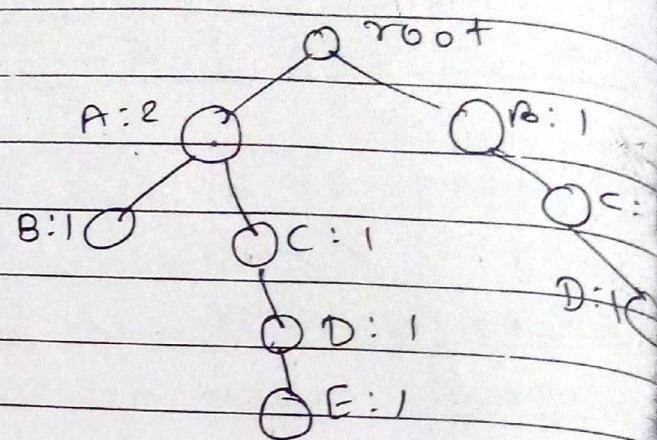
$$TID = 1$$



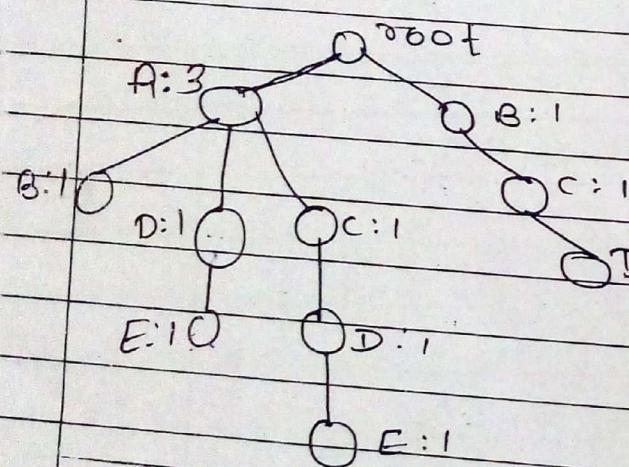
2. After reading TID = 2



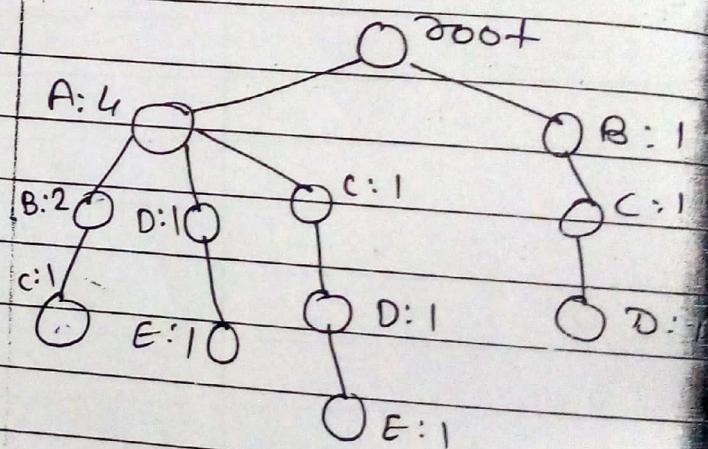
3. After reading TID = 3



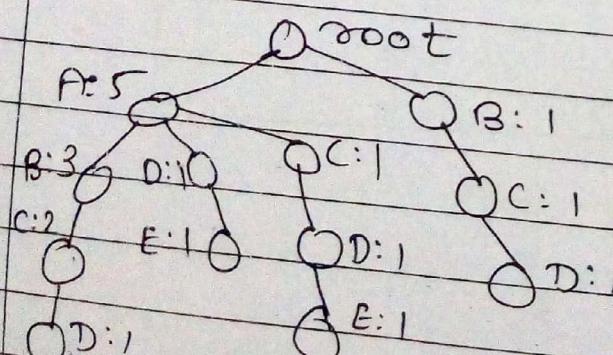
4. After reading TID = 4



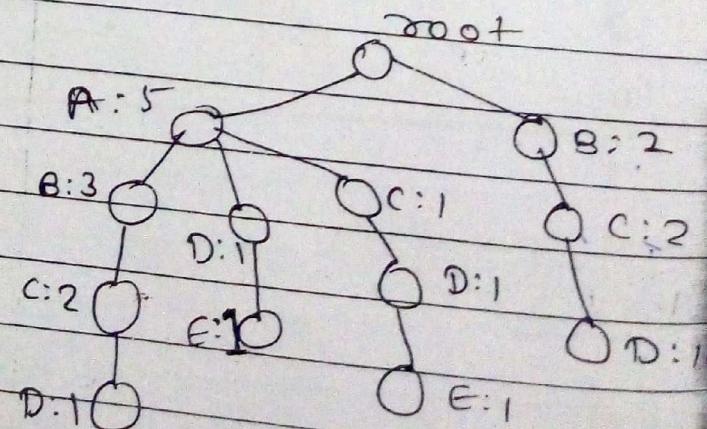
5. After reading TID = 5



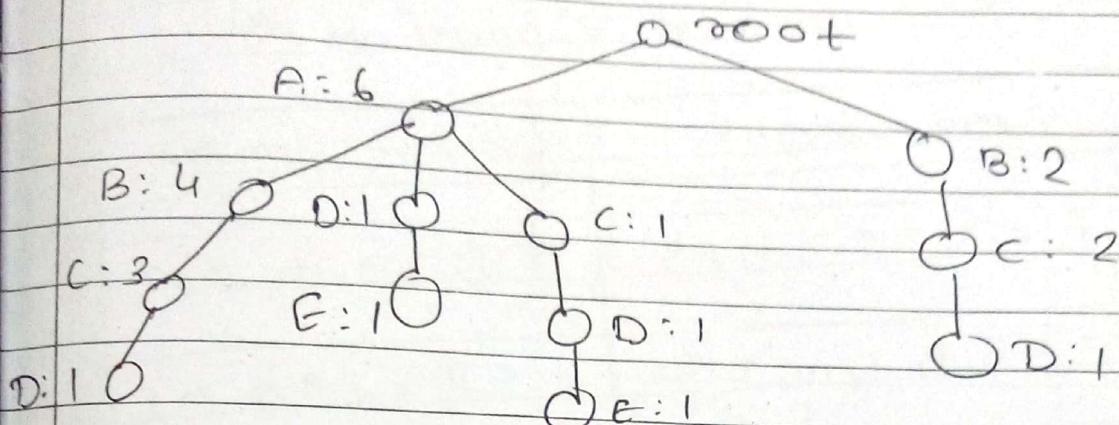
6. After reading TID = 6



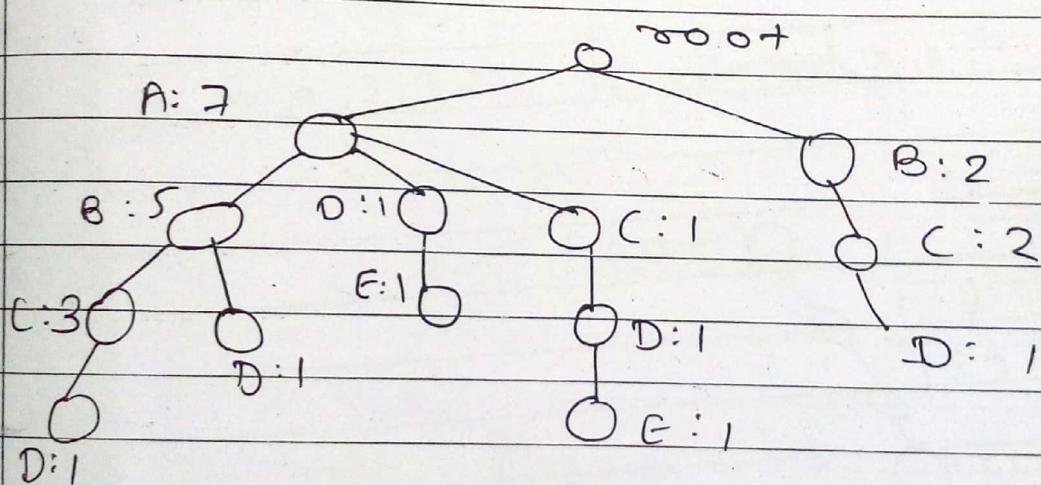
7. After reading TID = 7



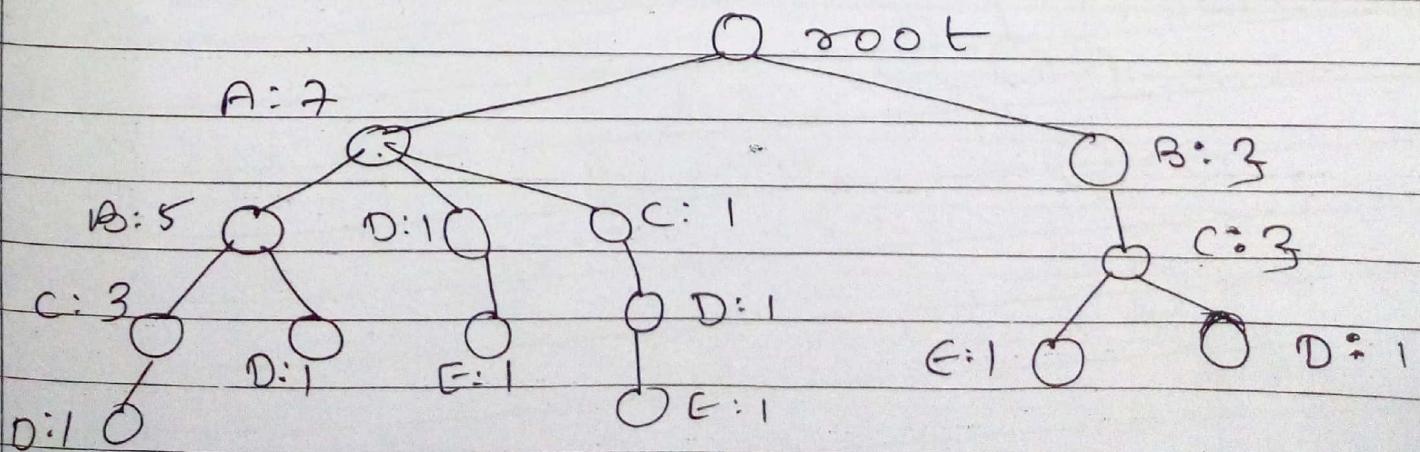
Q. After reading TID = 8



Q. After reading TID = 9



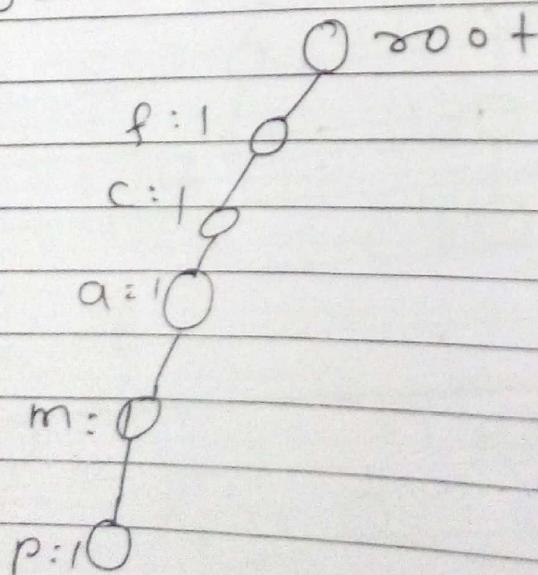
Q. After reading TID = 10



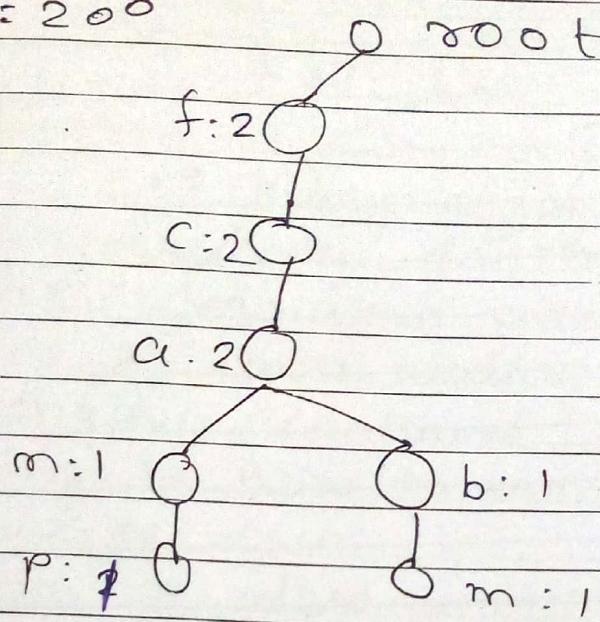
3. Construct FP-Tree for following transaction database with minimum support as one.

TID	Item bought	Frequent item set
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, n, i, o}	{f, b}
400	{b, e, k, s, p}	{f, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

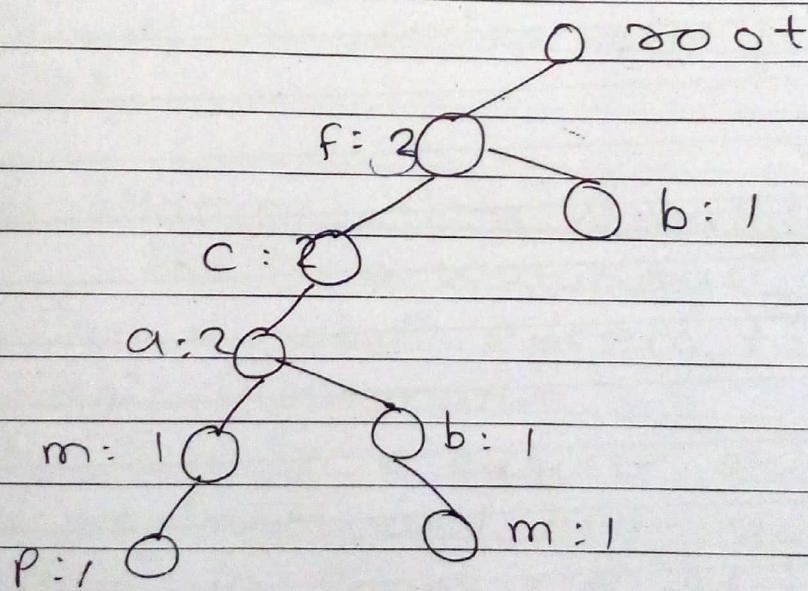
TID = 100



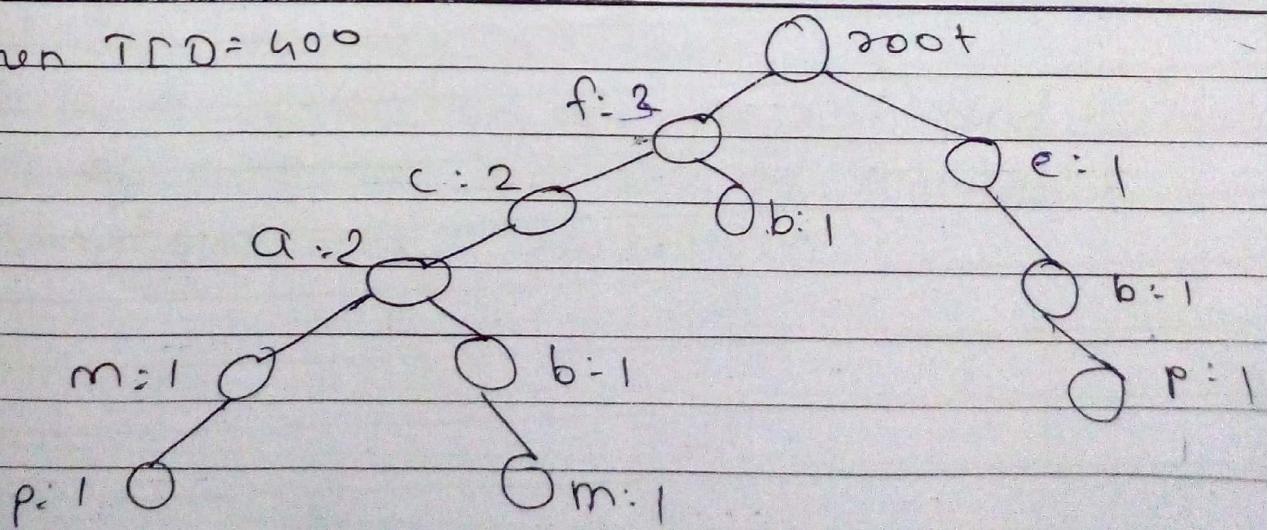
i) TID = 200



ii) TID = 300



when TID = 400



v) when TID = 500

