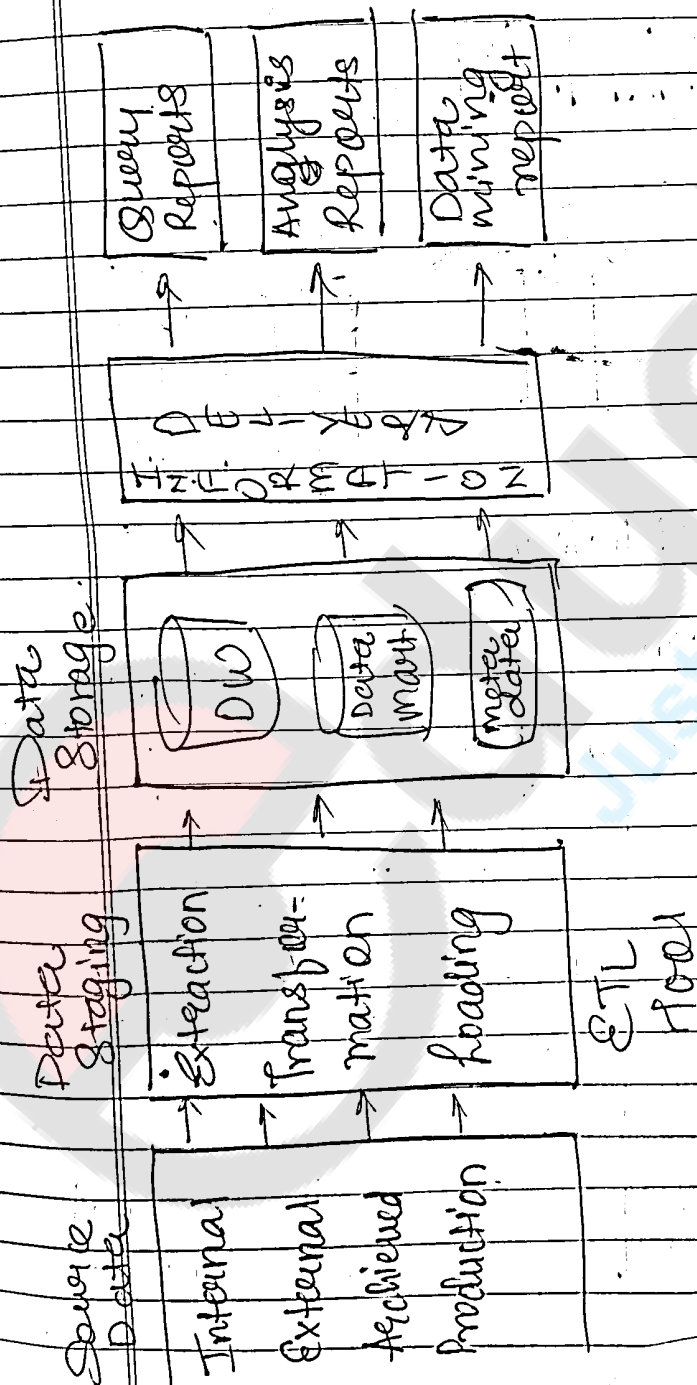


31<sup>st</sup> Jan

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

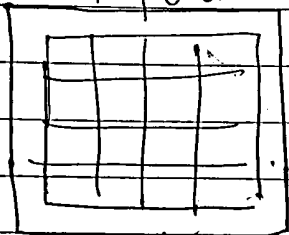
Data warehouse is a collection of subject oriented, integrated, non-volatile, time variant data.

\* DW architecture:

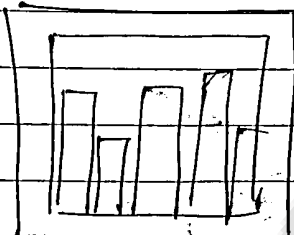


# \* 3-Tier Architecture of Data Warehouse.

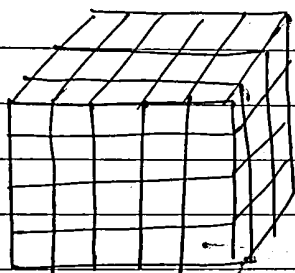
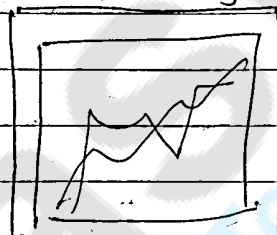
Query/  
Report



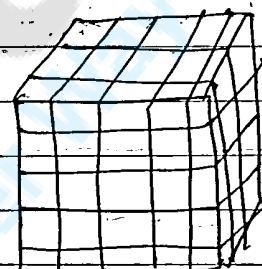
analysis



Data  
Mining



OLAP Server



OLAP Server

Middle Tier

Monitoring



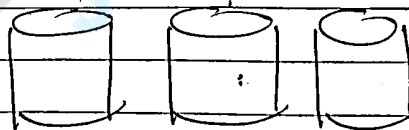
Administration



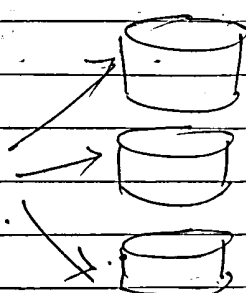
Metadata  
Repository



Data Warehouse

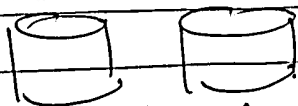
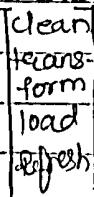


Data  
Marts

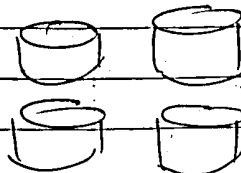


Data Warehouse Server

Extract Data

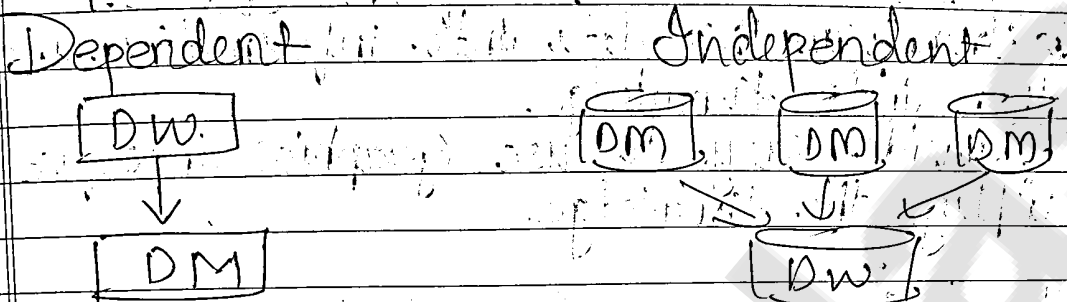


operational  
DB



External  
Source

## Data Mart

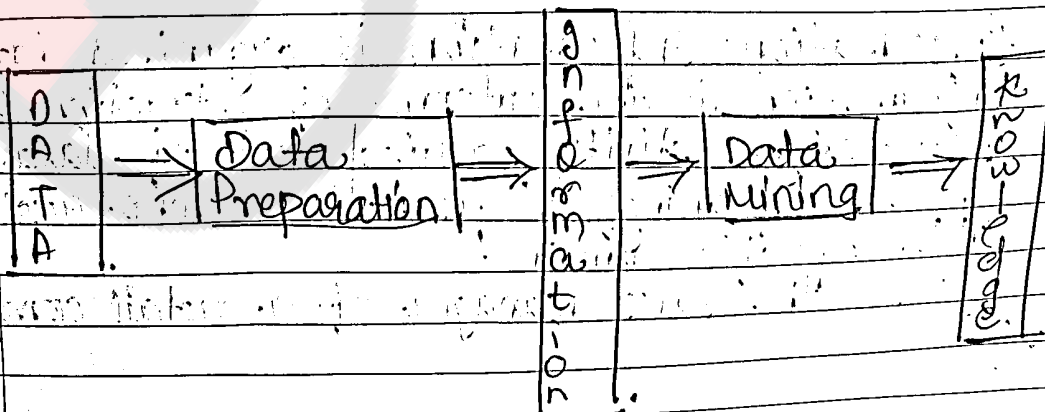


## Module 1 Business Intelligence

### \* Business Intelligence

BI is defined as a set of mathematical models and analysis methodology that exploit the available data to generate information & knowledge which is useful for complex decision making process.

- Q. Explain architecture of BI.  
 Q. Difference b/w Data, Information, Knowledge.



## \* Goal of BI system:

- ① Access data from different sources.
- ② Transform this data into information and then into knowledge.
- ③ Provide easy to use Graphics interface to display the knowledge

The main purpose of BI is to provide knowledge workers with tools and methodologies that allow them to make effective and timely decision.

## \* Difference b/w data, information & knowledge.

Data: Data are collected on a daily basis in the form of bits, number, symbols, objects, etc. Data represents a structural codification of single primary entity.  
Eg: for retail company primary entity or data are customer, point of sale (POS), salesman, product & sales receipt.

Information: Information is organized data. Information is the outcome of extraction & processing activity carried out of data & it appears meaningful for those who receive it in a specific domain.  
Eg: To the sales manager of a retail company

- 1) The proportion of sale of x product is 1000 per week.
- 2) The number of customer holding loyalty card is reduced to 50%.

### Knowledge:

Information is transformed into knowledge when it is used to make decision & develop the corresponding action.

Eg: Problem → For a Retail company, it is observed that they have reduced their usual amount of business in certain area.

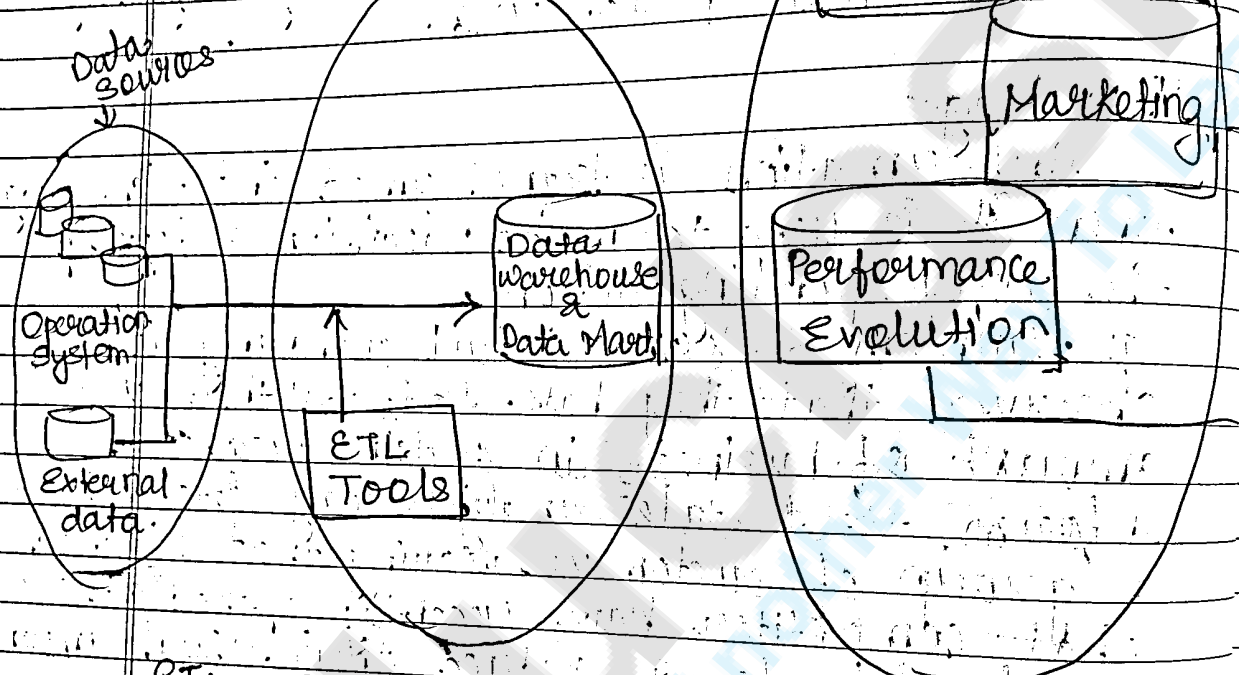
Infer → A sales analysis may detect that a group of customers, leaving that area is that competitor have recently opened a new point of sale & reduce their usual amount of sale.

Knowledge → The problem can be solved by introducing new services for the customer residing in that specific area.

Passive way → Knowledge Management  
Active way → BI

\* Imp  
10mks \*

BI architecture



ETL = Extraction Transformation & Load

Components of BI architecture

1] Data sources

In first stage it is necessary to gather & integrate the data stored in the various primary & secondary sources. Basically the data sources are operational sim, external data such as email, files, document etc.

2] Data Warehouse -

Using ETL tool, the data originating from different sources are stored in databases.

Q. What do you mean by BI & Architecture of BI.

- ↓
- 1) Multi Dimensional Cubes
  - 2) Exploratory data analysis
  - 3) Time Series analysis
  - 4) Data Mining
  - 5) Optimization
- ↑

which is referred as data-warehouses & data marts.

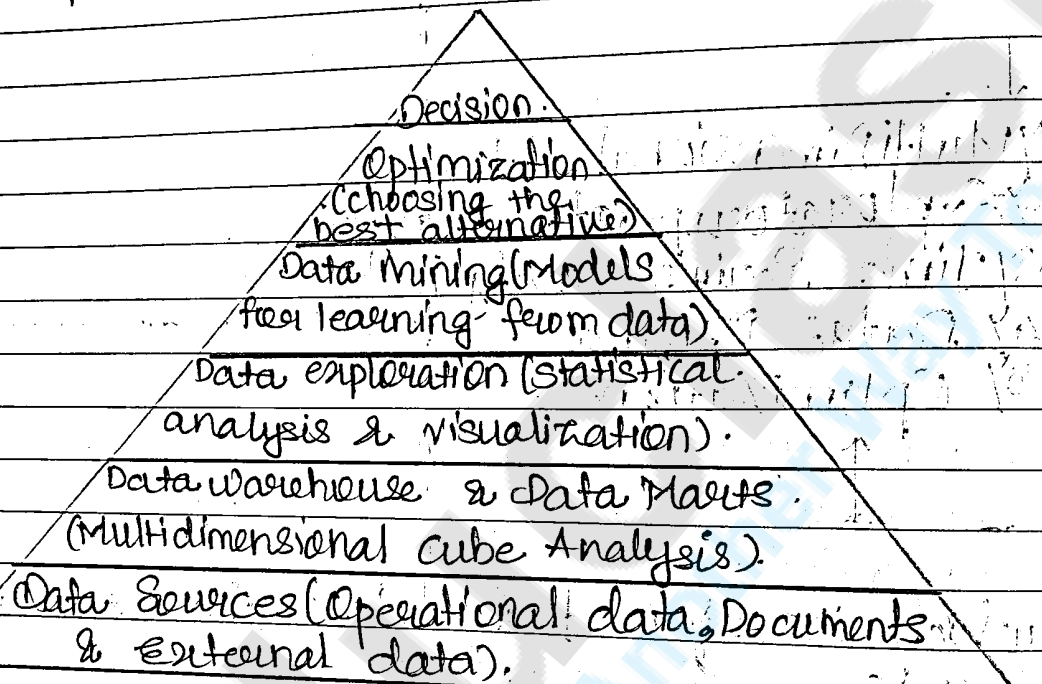
3] Business Intelligence:

Data are finally extracted and use to feed mathematical models & analysis methodologies to support decision.

Different decision support applications are:

- a) Multidimensional cubes,
- b) Exploratory data analysis,
- c) Time series analysis
- d) Inductive learning models for Data Mining.
- e) Optimization models

## \* Components of BI

Top to Bottom:

- Analyst & Experts
- ① Decision
  - ② Optimization
  - ③ Data Mining - Active BI Analysis
  - ④ Data Exploration - Passive BI Analysis
  - ⑤ DW & Data Marts
  - ⑥ Data Sources
- Decision Makers
- Db Administrators



Unit 8

(20m)

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## Introduction to DW

\* Differences b/w operational system & information system.

Operation System OR OLTP  
(Online Transaction  
Processing).

Information System OR  
OLAP (Online Analytical  
Processing).

1) Contains only current  
data values.

1) Contains archived,  
derived summarized data.

2) OS is optimized for  
transactions.

2) Optimized for complex  
queries.

3) Access frequency is  
high.

3) Access frequency is  
medium to low.

4) The data can be read,  
write, delete & update.

4) The data can only be  
read.

5) Large no. of users.

5) Very few no. of users.

Q. What do you mean by data warehouse & architecture of DW.

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## \* Define Datawarehouse.

Bill Inmon → 1990 → USA

↳ Father of Datawarehouse.

Datawarehouse is a subject oriented, integrated, non volatile, and time variant collection of data in support of management decision.

**Subject oriented data:** In the DW the data is stored by business objects or topics, not by operations applications like cust no, cust name address etc.

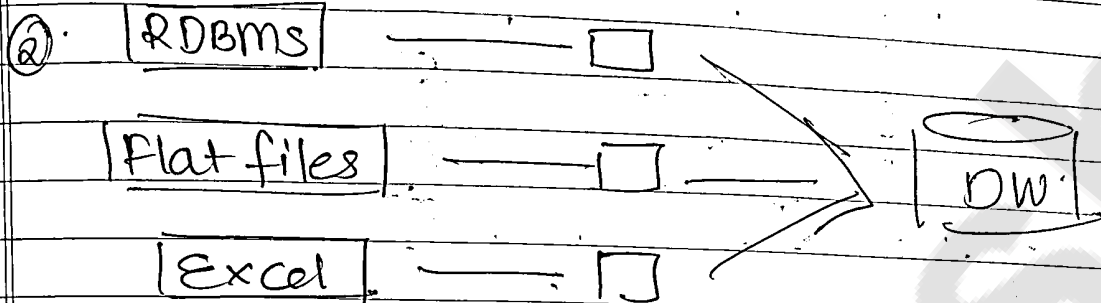
Eg: for manufacturing company, the subjects are customer, account, loan.

**Integrated data:** Integrated means the data are stored as a single unit, not as a collection of files that may have different structure.

**Non-volatile:** Non-volatile means the data do not keep changing in DW, once it is inserted in DW, it cannot be deleted.

① cust1      cust2      → [customer]

Subject oriented data.



## Integrated data

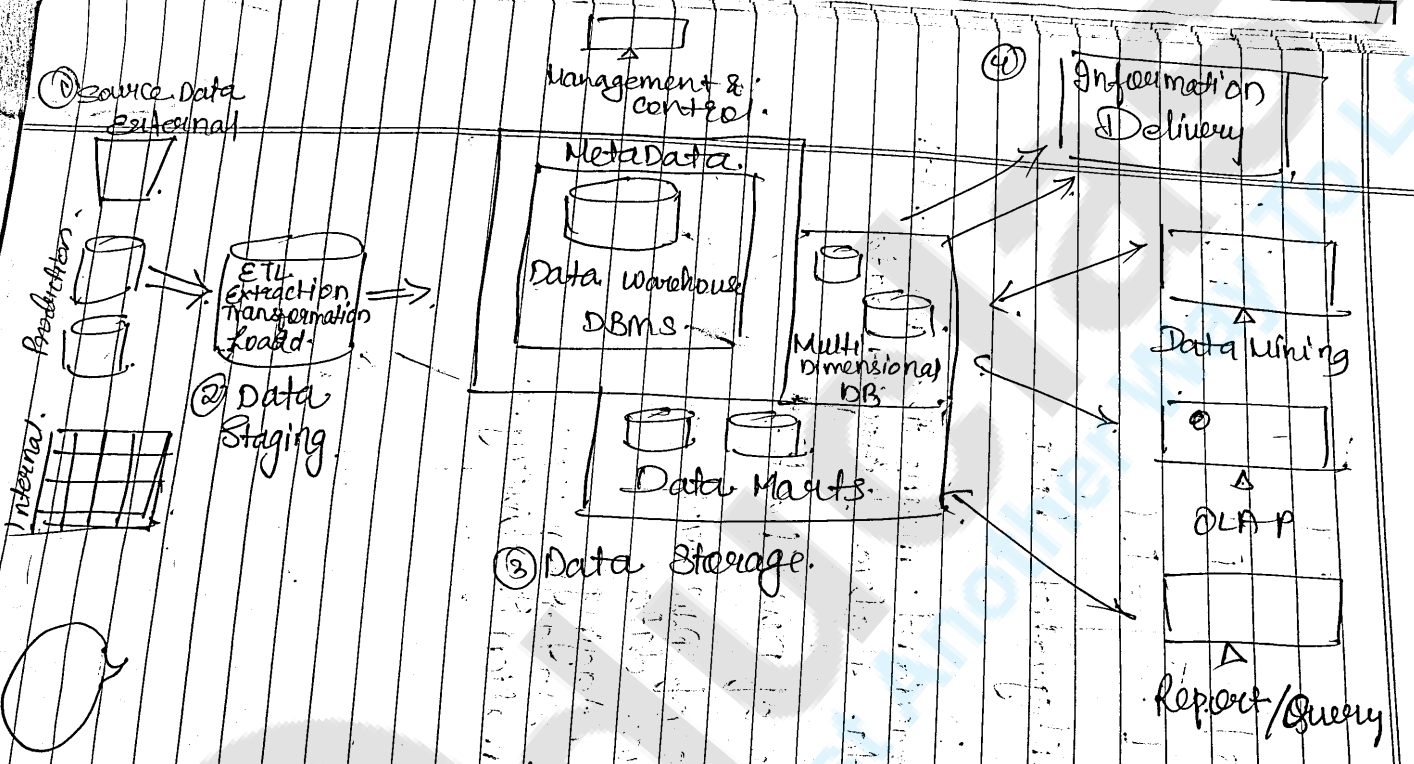
Time variant:

Time variant means the DW contains current data & also historical data.

\* Architecture of DW (back):

Components of DW are:

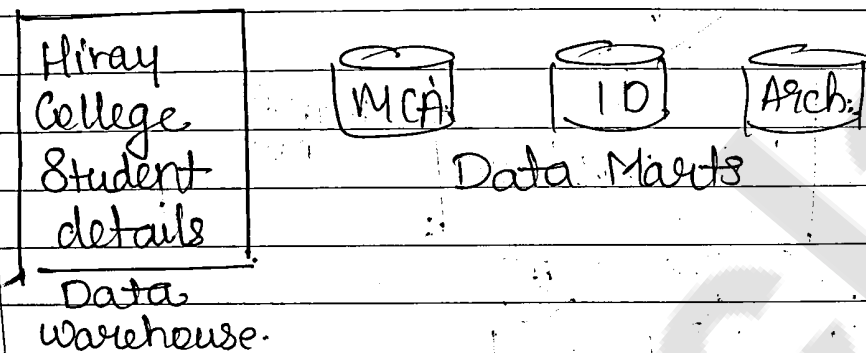
- ① Source data
- ② Data Staging
- ③ Data Storage
- ④ Information Delivery



DW components or Building Stages

- Q What is Data Warehouse? Explain the architecture.
- Q What is ETL? Explain it in detail.
- Q What is Metadata? Explain its types.
- Q What is Data Mart? Explain top down approach and Bottom up approach.

### SN. \* Data Mart (Departmental wise)



A subset of data warehouse that supports the requirements of a particular department or DW.

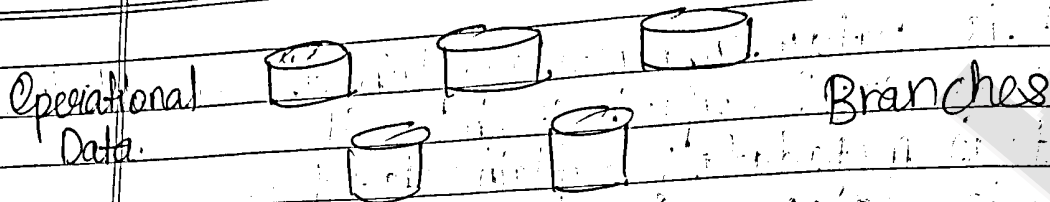
#### Characteristics of Data Mart

- ① It doesn't contain detail data like DW.
- ② More easily understood & navigated.
- ③ Can be dependent or independent.

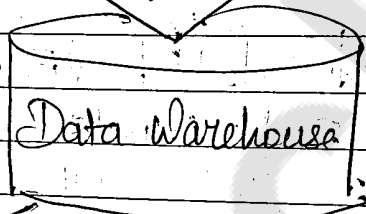
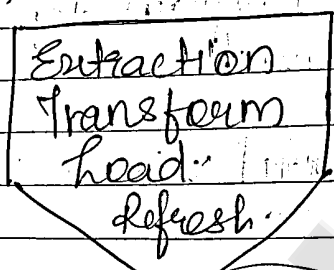
To develop data mart there are 2 approaches:

- ① Top Down Approach.
- ② Bottom up approach.

i) Top Down approach (Dependent Data Mart)

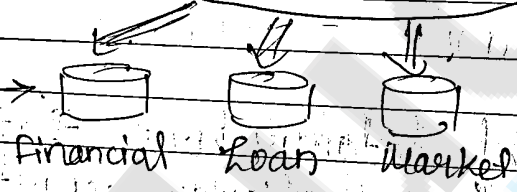


8/11



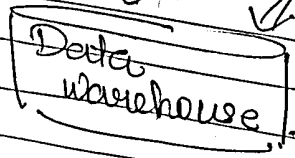
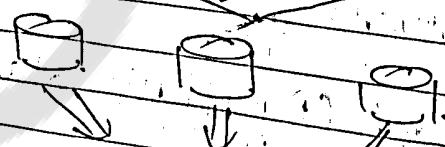
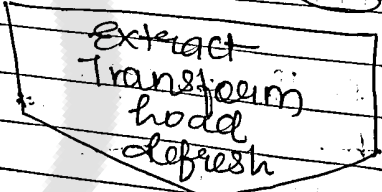
Bank DW

Data mart



2] Bottom Up Approach  
(Independent Data Mart)

Operational data



## 8/11/18 \* Metadata.

Metadata is a DW contains the answer to the questions about the data in the DW. The answer will be stored in a place which is called as meta data repository.

### Definitions of Metadata.

- ① Metadata is data about data.
- ② Metadata is table of contents for data.
- ③ Metadata is the DW road map.
- ④ Metadata is the glue that holds the DW content together.

Eg: Product Metadata // Metadata format.

Source: Finished goods orders,  
System Maintenance contracts,  
Online Sales.

Create date: Jan 15 2016

Last update date: Weekly

Last full refresh: Dec 29 2016  
data.

Full refresh cycle: Every year.

Data quality reviewed: Jan 25 2017.

Last Duplication: Jan 10 2018

Planned Archival: Every six months.

Responsible user: Jane Brown

## \* Types of Metadata

- ① Human Metadata
- ② Computer Based Metadata for user to use.
- ③ Computer Based Metadata for computer to use.

A metadata repository should contain the following information:

- A description of datawarehouse structure.
  - > Schema view
  - > dimension view
  - > data mart location & contents

- Operational Metadata
  - > history of migrated data
  - > Currency of data
  - > Monitoring Information

- Algorithm used for summarization
  - > Data granularity
  - > Partitions
  - > Aggregation

- Mapping from operational environment to DW
  - > Source Database
  - > gateway
  - > extraction
  - > cleaning



- Data related to system performance.
  - > How to improve data access.
  - > Rules for timing and scheduling to refresh
  - > updates.

- Business Metadata.
  - > Data ownership
  - > Business terms & definition.

### \* ETL Process.

10 MKS.

10/10/17

#### Extraction:

Extraction refers to the process of retrieving the required data from the operational system tables.

#### Transformation:

Transformation is the process which transforms the extracted data in accordance with the business rules & standards that have been established for the data warehouse.

Transformation techniques are:

1) Format changes.

DD mm yy.

yy mm DD.

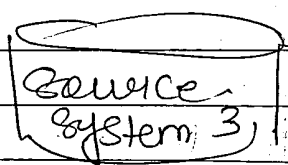
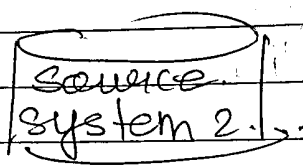
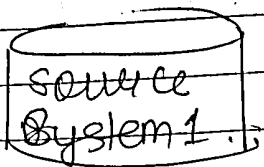
2) Duplications.

3) Splitting up fields.

4) Integrating fields.

5) Replacement of values.

Relational Database.



File Sources.

- Excel (.xls) files
- Text (.txt) files
- XML (.xml) files
- Other files

Other Sources

D  
A  
T  
A  
f  
r  
o  
m  
D  
i  
f  
f  
e  
r  
e  
n  
t  
S  
o  
u  
r  
c  
e  
s

Rec  
Ge

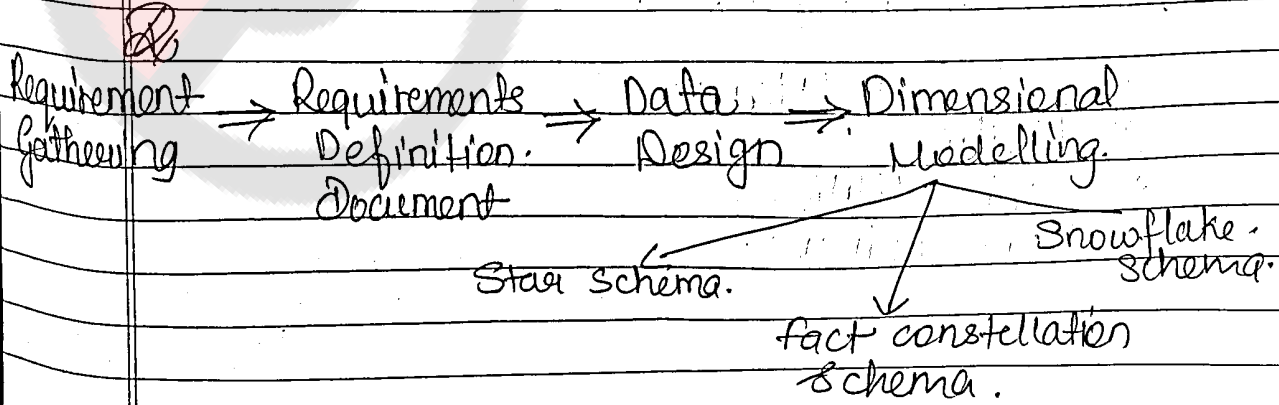
loading:

X Dimension Modeling

Q Difference between ER modeling & dimensional modeling.

ER Modeling	Dimensional Modeling
① Remove data redundancy	① Captures critical measures
② Ensures data consistency	② Views data along dimensions
③ Express microscopic relationship	③ Intuitive to business users.

Q what do you mean by dimensional modelling with eg: [with 3 schemas]



- Dimensional modeling is a designing concept used by many datawarehouse designers to build DW.

- In dimensional modeling two types of tables are used:

> Fact table

> Dimension table

Fact table:

- Fact tables are used to record actual facts or measurement units in the business with no redundancy.

- Facts are numeric data item that are of interest to the business.

- Eg: In retail: No. of units sold, sales amount  
telecommunication: length of call in minutes, average no. of calls.

Banking: average daily balance, transaction amount.

Dimension table:

1. Dimension table establish the context of the facts.
2. It stores the fields that describe the facts.

Types of facts (Measure)

① Additive.

② Semi additive.

③ Non additive.

Star schema design →

Fk + measurement  
Should be normalized

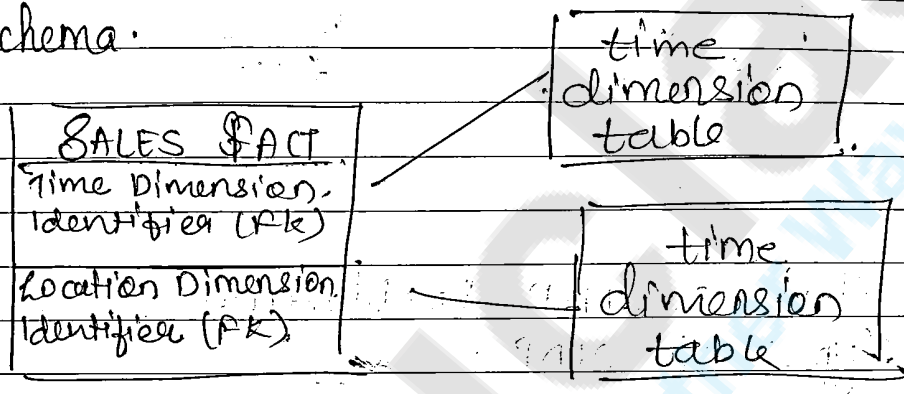
Fact table connected with multiple dimension.

### 8 Telecommunication system:

- ① fact table
- ② Dimension table

### \* Star Schema

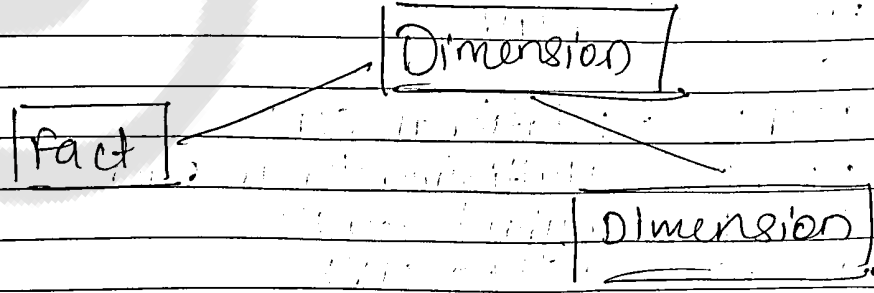
Smks



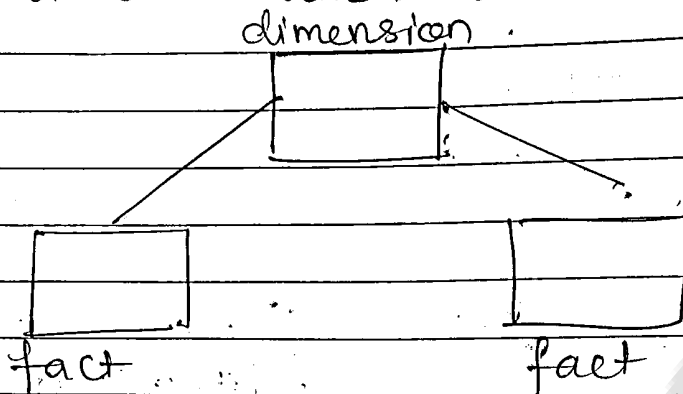
- 1) It is one type of dimensional model.
- 2) Used to view cube data.
- 3) One fact table connected with multiple dimension.
- 4) Dimension table may not be normalized.

### A Snowflake Schema

— Extension of star schema.



## \* Fact Constellation Schema



Q What is OLAP Architecture

Q Types of OLAP

\* OLAP - Online Analytical Processing

- OLAP is a SW frontend tool for datawarehouse  
 - It is an information delivery system for DW  
 - OLAP SW provides the ability to analyze the large volume of information to improve decision making at all level of an organization.

\* Types of OLAP

ROLAP → Relational OLAP

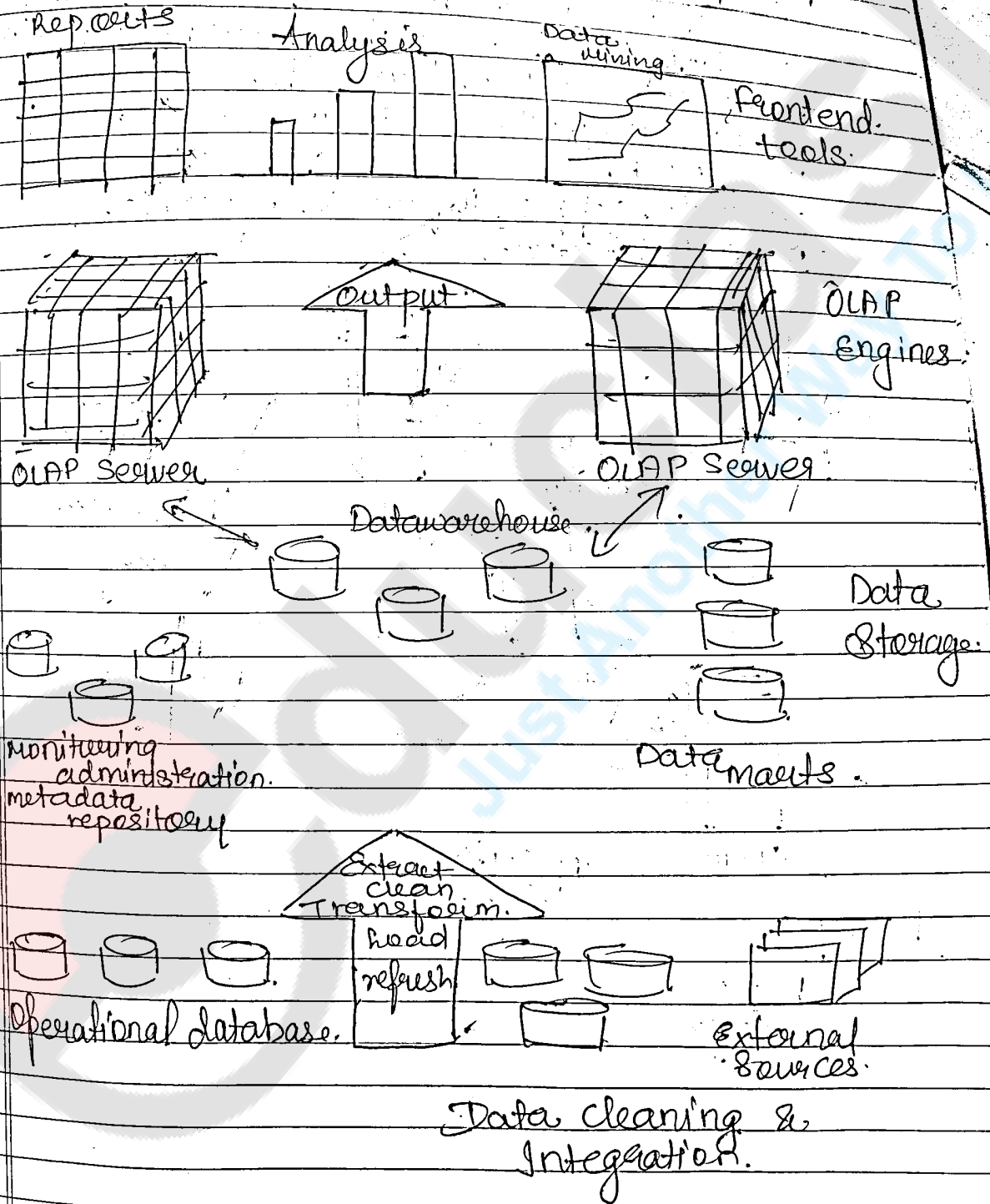
MOLAP → Multidimensional OLAP

HOLAP → Hybrid OLAP

DO LAP → Desktop OLAP

WOLAP → Webbased OLAP

# \* OLAP Architecture.

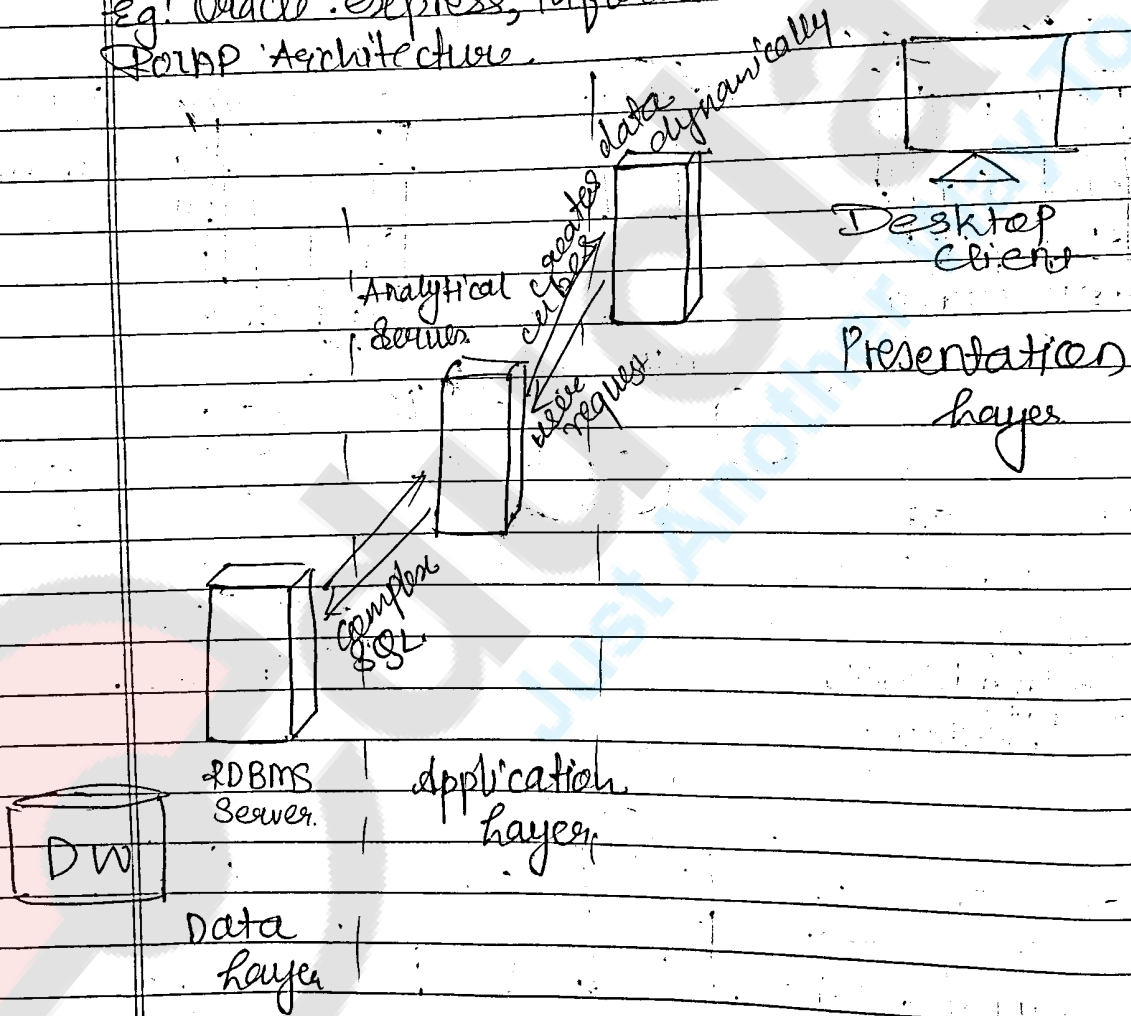


### 17) ROLAP

ROLAP is a s/w which is used to retrieve the data from a database-oriented relationship. It is suitable for large db. ROLAP has ability to analyse a business better than MOLAP.

Eg: Oracle express, Informatica metacube.

### ROLAP Architecture

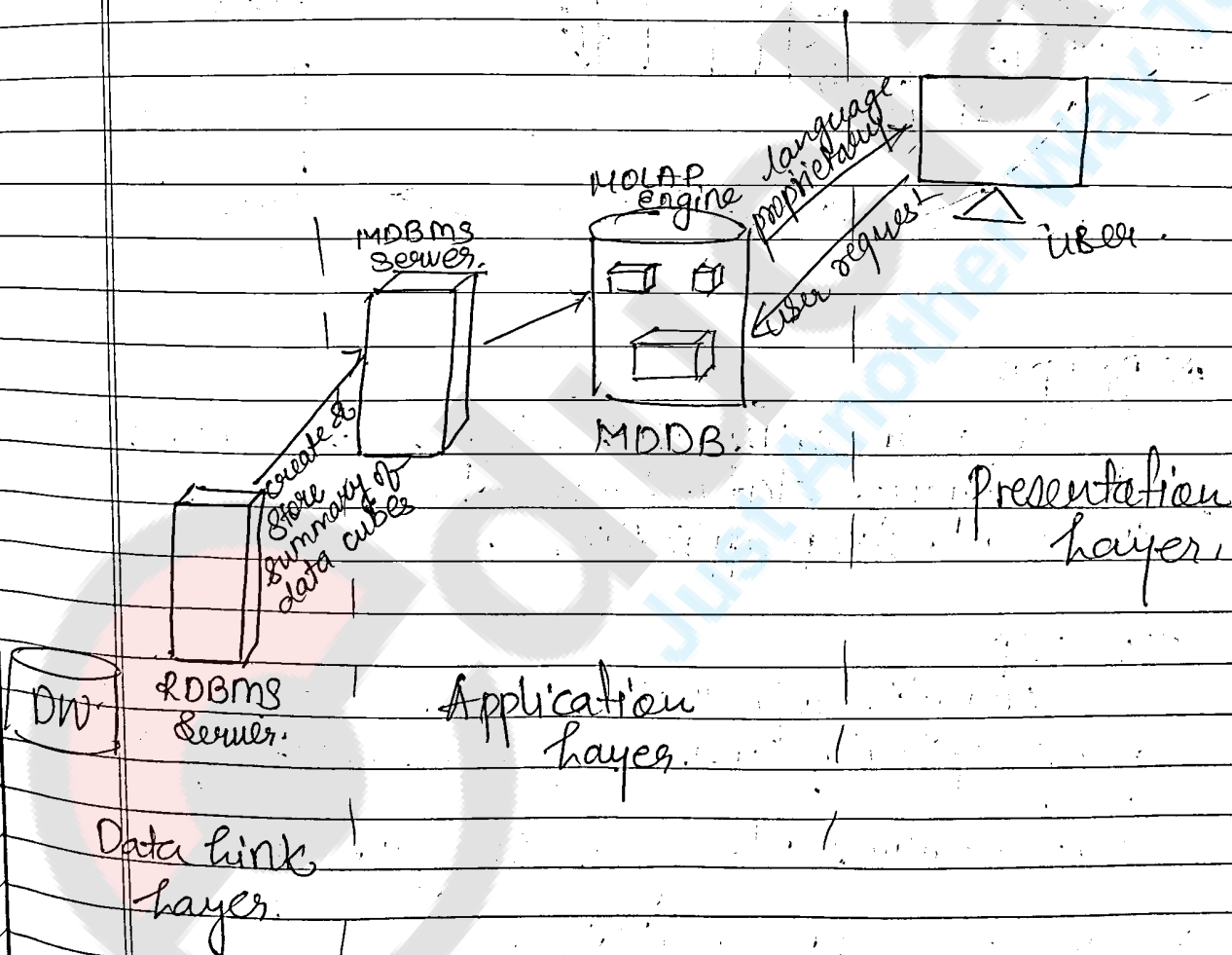




Q7 MOLAP

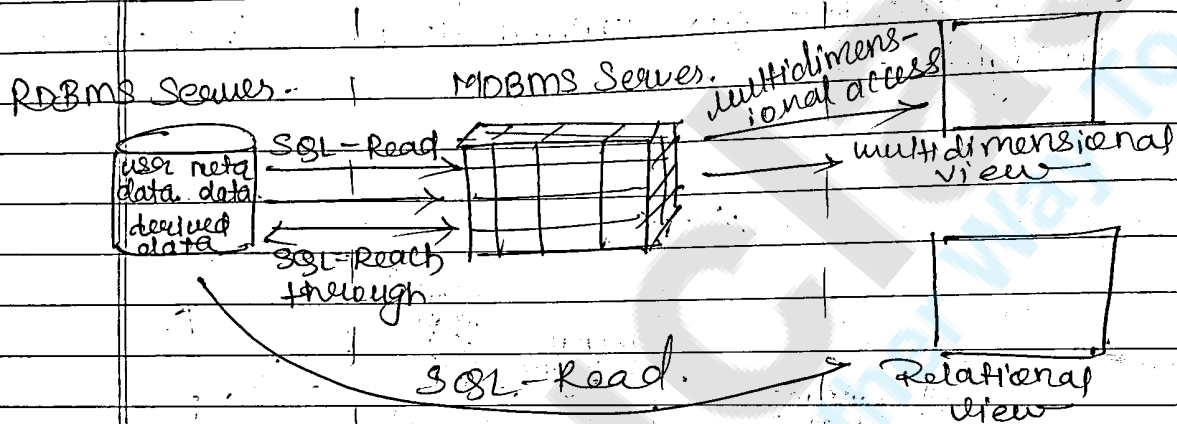
It is a s/w tool which is used to retrieve multidimensional data cubes from database. It is faster than ROLAP but it is suitable for small database.

Eg: Microsoft s/w commander, FDC, Dimensional insight etc.



3) HOLAP

- used to retrieve both relational & multi-dimensional relation data from the db.  
Eg: SAS, CFS, Vision, Arbor, Essbase.

4) DOLAP

used in only one single standalone pc to retrieve data from db.  
It is a desktop application.

5) WOLAP

Analytic information is embedded into html pages, specially written for the needs of special users.

Eg: Executives, Managers, Analysts, Specialists

\* Types of OLAP operations:

- Roll-up (drill-up)

- \* - Drill-down (roll down)
- Slice & Dice
- Pivot (rotate)

# mfb

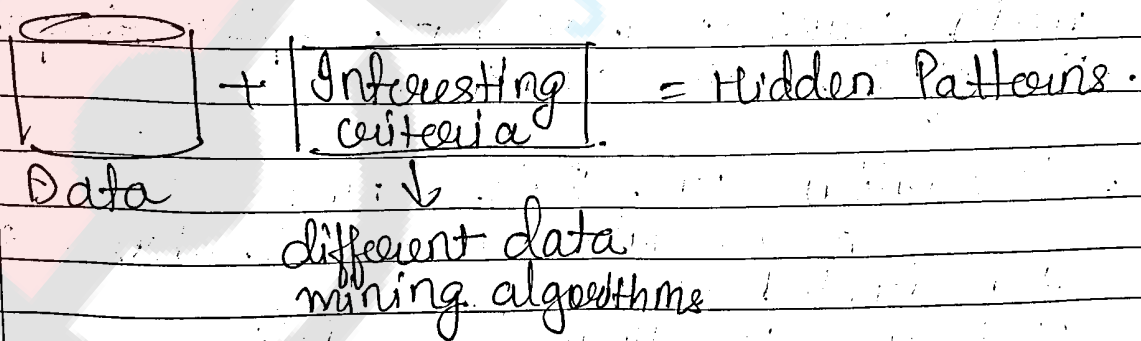
\* Data Mining:

- DM is a technique to find the hidden data from dw.

- DM is a collection of techniques for efficient automated discovery of previously unknown valid and understandable patterns in large databases.

The patterns must be actionable so that, they can be used for decision making process.

- DM is the generic term used to look for hidden patterns in db.



It has been also called exploratory data analyses, data driven discovery & deductive learning process.

## Data Mining Model:

### Predictive

- Classification
- Regression
- Time Series
- Prediction

### Descriptive

- Clustering
- Summarization
- Association Rules
- Sequence Discovery

Predictive model makes a prediction about values of data using known results found from different data predicting modelling or maybe from the other historical data.

Descriptive model identifies patterns or relationship in data. Unlike the predictive model, a descriptive model serves as a way to explore the properties of the data examined.

### x Association Rule Data Mining

An association is a rule when  $x$  then  $y$ . It is denoted as  $x \rightarrow y$ . It is used for finding the relationship among the data.

eg: If India wins in cricket world cup then sales of sweets rises.

Suppose  $x = \text{India wins Cricket}$ .  
 $y = \text{Sales of sweets rises}$   
 $x \rightarrow y$ .

In association rule mining we use two terms:  
 1) Support 2) Confidence.

1) Support: probability that a transaction contains  $x$  &  $y$  i.e. it is the percentage of transaction in which item  $x$  &  $y$  occurred together.

2) Confidence: it is a conditional probability that transaction having  $x$  also contains  $y$ .

Confidence is a rule for finding the ratio of the no. of occurrence of  $y$  given  $x$  among all other.

Q Find the support & confidence for the rule  $\text{Bag} \rightarrow \text{uniform}$  by using the full database.

$t_1 \rightarrow \text{Bag, uniform, Crayons}$

$t_2 \rightarrow \text{Pencil, uniform, Bag, books}$

$t_3 \rightarrow \text{Book, Bag, uniform}$

$t_4 \rightarrow \text{Bag, Pencil, Book}$

$t_5 \rightarrow \text{Bag, uniform, pencil}$

$t_6 \rightarrow \text{Uniform, crayons, Bag}$

$t_7 \rightarrow \text{Bag, Pencil, Book}$

$t_8 \rightarrow \text{Crayons, uniform, Bag}$

$t_9 \rightarrow \text{Books, crayons, Bag}$

$t_{10} \rightarrow \text{Uniform, crayons, pencil}$

$$\text{Support} = \frac{\text{No. of transactions containing Bag \& uniform}}{\text{total no. of transaction}} \\ = \frac{5}{10} \times 100 = 50\%$$

$$\text{Confidence} = \frac{\text{No. of transactions both bag \& uniform}}{\text{No. of transactions containing bag}} \\ = \frac{5}{8} \times 100 = 62.5\%$$

Q find support & confidence for the rule  
Milk  $\rightarrow$  Butter

t1  $\rightarrow$  Milk, Bread, Butter

t2  $\rightarrow$  Milk, Butter

t3  $\rightarrow$  Jam, Butter

t4  $\rightarrow$  Bread, Egg, Fruit

$$\text{Support} = \frac{2}{4} \times 100$$

$$= 50\%$$

$$\text{Confidence} = \frac{2}{2} \times 100$$

$$= 100\%$$

Q Confiden t<sub>1</sub> = {F, A, D, B}

t<sub>2</sub> = {D, A, C, E, B}

t<sub>3</sub> = {C, A, B, E}

t<sub>4</sub> = {B, A, D}

{B, D}  $\rightarrow$  A

$$\text{Support} = \frac{3}{4} \times 100$$

$$= 75\%$$

$$\text{Confidence} = \frac{3}{3} \times 100$$

$$= 100\%$$

28/02

Q What is the need of data preprocessing?  
Techniques of preprocessing.

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

Data Quality:

- Measures for data quality: a multidimensional view.
- accuracy.
- completeness.
- consistency.
- timeliness.
- believability.
- interpretability.

Techniques of Preprocessing:

- Data cleaning
- Data integration
- Data reduction
- Data transformation and data discretization.

① Data cleaning:

Data in real world is dirty: lots of potentially incorrect data. eg. instrument fault, human or computer error.

- incomplete
- noisy.
- inconsistent
- intentional.

\* How to handle missing (incomplete) data:

- Data is not always available.
- Missing data may need to be inferred.
- Missing data may be due to:

- > equipment malfunction.
- > data not entered due to misunderstanding
- > inconsistent with other recorded data & thus deleted.
- > not register history or changes of the data.

- Ignore the tuple.

- Fill in the missing value manually.

- Fill in it automatically with:

> global constant.

> attribute mean.

> the most probable value: reference-based Bayesian formula or decision tree.

> attribute mean for all samples belonging to same class.

## Noisy Data.

Noise: random errors or variance in a measured variable.

Incorrect attribute values.

> faulty data collection instruments

> data entry

> data transmission

> technology limitation

> inconsistency in naming convention

Other data problems → requires data cleaning



How to handle noisy data.

- Binning
- Regression
- Clustering
- Combined computer & human inspection.

## Binning

- Binning method smooths sorted data values by consulting its neighbourhood i.e. the values around it.

Method:

- ① Sort the data in descending or ascending order.
- ② Partition the sorted data into equal frequency bins
- ③ Apply any binning method such as - smooth by bin means,
  - Smooth by bin median
  - Smooth by bin boundaries, etc.

Q. Suppose a group of 9 sales price records has been sorted as follows.

4, 8, 15, 21, 21, 24, 28, 28, 34.

By considering bin depth of 3 use smooth by bin means and smooth by boundaries.

→ Partition the given data in equal freq bins i.e 3.

Bin 1: 4, 8, 15

Bin 2: 21, 21, 24

Bin 3: 28, 28, 34.

1) Smooth by <sup>bin</sup> mean.  
- Find mean of each <sup>bin</sup> ~~mean~~ & replace the values of each bin.

Bin 1:  $\frac{4+8+15}{3} = \frac{27}{3} = 9$

Bin 2:  $\frac{21+21+24}{3} = 22$

Bin 3:  $\frac{25+28+34}{3} = 29$

After smoothening the bin values are:

Bin 1: 9, 9, 9

Bin 2: 22, 22, 22

Bin 3: 29, 29, 29

ii) Smooth by bin boundaries  
Find the max & min of each bin.

Bin 1: 4, 8, 15      min = 4      max = 15

2: 21, 21, 24      21      24

3: 25, 28, 34      25      34

Replace the data <sup>in bin</sup> by boundary values

Bin 1: 4, 4, 15

2: 21, 21, 24

3: 25, 25, 34

1) apply smooth by bin means & smooth by bin boundaries, smooth the foll data.  
Bin depth = 4.

10, 4, 3, 10, 11, 6, 7, 5, 6, 12, 15, 16

→ 1) Sort the data:

3, 4, 5, 6, 6, 7, 10, 10, 11, 12, 15, 16.

2) Partition

Bin 1: 3, 4, 5, 6

Bin 2: 6, 7, 10, 10

Bin 3: 11, 12, 15, 16

3) Smooth by bin mean

$$\text{Bin 1: } \frac{3+4+5+6}{4} = 4.5$$

$$\text{Bin 2: } \frac{6+7+10+10}{4} = 8.25$$

$$\text{Bin 3: } \frac{11+12+15+16}{4} = 13.5$$

after smoothening the bin values are:

$$\text{Bin 1: } 4.5, 4.5, 4.5, 4.5$$

$$\text{2: } 8.25, 8.25, 8.25, 8.25$$

$$\text{3: } 13.5, 13.5, 13.5, 13.5$$

4) Smooth by bin boundary

Bin 1: 3, 4, 5, 6

min = 3

max = 6

2: 6, 7, 10, 10

min = 6

max = 10

3: 11, 12, 15, 16

min = 11

max = 16

Scan

After Smoothing:

- Bin 1: 3, 3, 6, 6
- 2: 6, 6, 10, 10
- 3: 11, 11, 16, 16

— Regression method is used to smooth by fitting the data in regression function.  
— Linear regression used to find the best line to fit two attributes so that one attribute can be used to predict the other.

Linear multiple regression: extension of linear regression where more than two attributes are involved & the data are fit to multidimensional surface.

— Clustering: <sup>Euclidian</sup> ~~Euclidian~~ By using ~~Euclidian~~ distance formula the distance between outliers & clusters will be found out then the nearest cluster will be considered for the outlier.

## \* Data Integration

Combines data from multiple source into a coherent store.

### Schema Integration

Integrate metadata from different source.

### Entity Identification problem:

- Identify real world entities from multiple data sources.

Detecting and resolving data value conflicts for the same real world entity, attribute values from different sources are different.

Possible reasons: different representations.

- are any attributes correlated
- Tuple duplication
- Value confliction

#

Apriori algorithm

is used to discover all the frequent item sets from a large db in an efficient manner. The idea behind Apriori algorithm is: if all item set is frequent then all subsets of this item set should also be frequent.

Q Find frequent item sets or association rule with 30% support & 70% confidence of foll db:

Transaction I.d.	Item Set
t1	Bag, uniform, crayons
t2	Books, Bag, uniform
t3	Bag, uniform, Pencil
t4	Bag, pencil, Books
t5	Uniform, crayons, Bag
t6	Bag, Pencil, Books
t7	Crayons, uniform, Bag
t8	Books, crayons, Bag
t9	Uniform, Crayons, pencil
t10	Pencil, uniform, Books

Soln:

Given support = 30% i.e out of 10 transaction atleast  $\frac{10}{100} \times 30 = 3$  or more than 3 frequency should be considered.

1<sup>st</sup> Item Set

Items	Frequency
Bag	8
Uniform	7
Crayons	5
Books	5
Pencil	5

Since all the items have more than 3 frequency they will be considered for the 2<sup>nd</sup> frequency item set.

2<sup>nd</sup> frequency Item Set

Items	Frequency
Bags, uniform	5
Bags, Crayon	4
Bags, Books	4
Bags, Pencil	3
Uniform, Crayons	4
Uniform, Books	2
Uniform, Pencil	3
Crayons, Books	1
Crayons, Pencil	1
Books, Pencil	3

Since {uniform, Books, Crayon, book} & {crayon, pencil} has frequency  $\leq 3$   $\therefore$  it would be discarded and others will proceed to the next level.

### 3<sup>rd</sup> frequency Item Set

Items	frequency
Bag, Uniform, Crayons	3
Bag, Uniform, Books	1
Bag, Uniform, pencil	1
Bag, Crayon, books	1
Bag, Crayons, pencil	0
Bag, books, pencil	2
Uniform, Crayons, pencil	1
Uniform, pencil, books	1

from the Item set only Bag, uniform, crayons satisfies the frequency 3. i.e support 30%. Other item sets will be discarded.

Since there is only one frequency item set left, there will be no further combination.

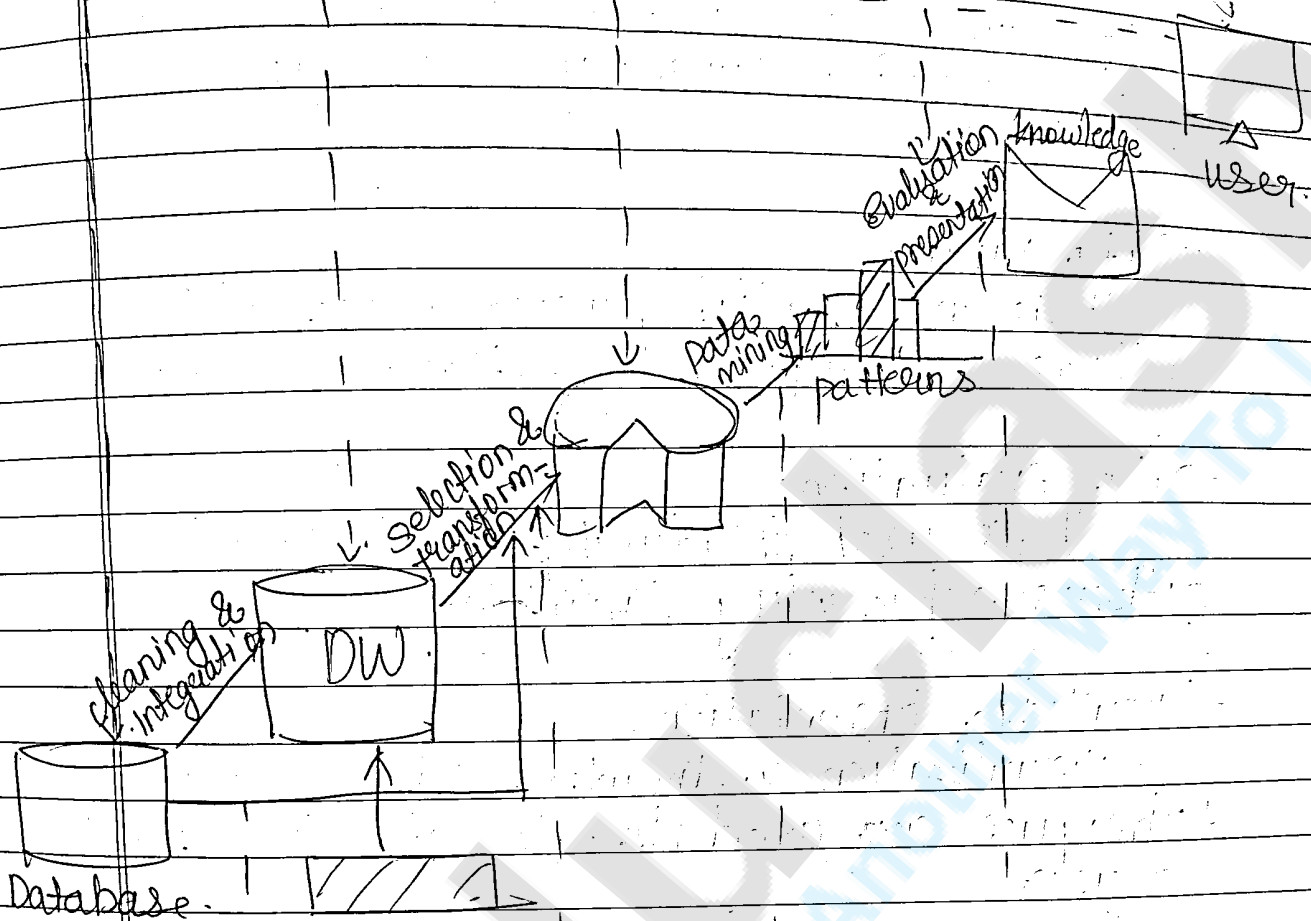
22/02  
VAMP

### KDD process (Knowledge discovery from database)

- KDD is the process of finding useful information & patterns in data

- Data Cleaning & Integration:  
Noisy data & inconsistent data will be removed from the db.





- Data Integration:

Multiple data sources are combined into a single form.

- Data Selection:

In this step the relevant data to the analysis task are retrieved from the db.

- Data Transformation:

In this step data are transformed & consolidated

into forms which are appropriate for mining while performing summary & aggregation appreciation.

- Data Mining:

It is an essential process where intelligent methods are applied to extract data patterns.

- Data evaluation:

It is used to identify the truly interesting patterns representing knowledge.

- Knowledge presentation:

Visualization & knowledge representation techniques are used to present mined knowledge to users.

1st march

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## Data Integration:

problems:

\* Handling redundancy in Data Integration.

- Correlation analysis is used for nominal data.

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} \quad \text{--- (one dimensional)}$$

correlation relationship between two attributes A & B can be discovered by  $\chi^2$  test.

Suppose A has c number of distinct values i.e.  $a_1, a_2, \dots, a_c$ .

and B has r distinct values represented as  $b_1, b_2, \dots, b_r$ .

To perform  $\chi^2$  test, a contingency table will be organised with c values of A making columns & r values of B making of the rows.

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad \text{--- (two dimensional)}$$

$$E_{ij} = \frac{\text{count}(A=a_i) \cdot \text{count}(B=b_j)}{n}$$

where  $n = \text{no. of tuples}$

$\text{count}(A=a_i) \Rightarrow \text{no. of tuples having value } a_i$

$\text{count}(B=b_j) \Rightarrow \text{no. of tuples having value } b_j$

The test is based on a significant level with  $(r-1) \times (c-1)$  degrees of freedom.  
If the hypothesis is rejected then we say that A & B statistically correlated.

Q Correlation analysis of nominal attribute using  $\chi^2$

Suppose a group of 1500 people were surveyed. The gender of each person was noted. Each person was polled as to whether his or her preferred type of reading material was fiction or non-fiction.

The following table gives the contingency table data for the surveyed data. Find are gender & preferred reading correlated?

	Male	Female	
fiction	250 (90)	200 (360)	450
non-fic	50 (120)	1000 (1340)	1050
	300	1200	1500

$$\rightarrow C_{11} = \text{count}(\text{male}) * \text{count}(\text{fiction})$$

$$= 250 * 300 + 450 = \frac{135000}{1500} = 90$$

Since our calculated value (507.82) > 10.828  
 $\therefore$  the hypothesis is rejected.  
It will be concluded that the 2 attributes are strongly correlated.

$$e_{12} = \frac{\text{count}(\text{female}) \times \text{count}(\text{fiction})}{n}$$

$$= \frac{1200 \times 450}{1500} = 360$$

$$e_{21} = \frac{\text{count}(\text{male}) \times \text{count}(\text{non-fic})}{n}$$

$$= \frac{300 \times 1050}{1500} = 210$$

$$e_{22} = \frac{\text{count}(\text{female}) \times \text{count}(\text{non-fic})}{n}$$

$$= \frac{1200 \times 1050}{1500} = 840$$

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840}$$

$$= \frac{25600}{90} + \frac{25600}{210} + \frac{25600}{360} + \frac{25600}{840}$$

$$= 284.44 + 121.90 + 71.11 + 30.47$$

$$= \underline{\underline{507.92}}$$

Degree of freedom:  $(r-1) \times (c-1)$

$$= (2-1) \times (2-1)$$

$$= 1 \times 1 = 1.$$

For 1 degree of freedom the  $\chi^2$  value needed to reject the hypothesis at 0.001 significance level is 10.828

class fruit & achieve polymorphism into it.

13/02

DMBI

Q Given a simple transactional database. Find FP tree for this database if minimum support threshold is 3.

TID	<del>Itemset</del> Itemset
01	f, a, c, d, g, i, m, p
02	a, b, c, f, h, m, o
03	b, f, h, j, o
04	b, c, k, s, p
05	a, f, c, e, d, p, m, n

Sol:

Scan the database & find frequent Item having frequency 3.

Items	frequency
f	4
a	3
c	4
d	1
g	1

1	Qualified items are
3	$\{(f, 4), (a, 3), (c, 4), (m, 3),$
3	$(p, 3), (b, 3)\}$
2	
2	
1	
1	
1	
1	
1	

- FP growth model.

Arrange the items in descending order of frequency.

$$L = \{(f, 4), (c, 4), (a, 3), (b, 3), (m, 3), (p, 3)\}$$

list of frequent items

\* Developing FP tree.

The root of the tree will be labelled as ROOT.

Scan the db second time.

Scan the first transaction which leads to the construction of first branch of FP tree.

Scan of first transaction  $(f, a, c, d, g, m, p)$   
only those items that are in the list of frequent items ~~only those are~~  
~~items~~ are selected.

i.e.  $\{(f,1), (c,1), (a,1), (m,1), (p,1)\}$

[Note: sequence as L is compulsory]

<sup>FP</sup>  
The tree will be.

ROOT

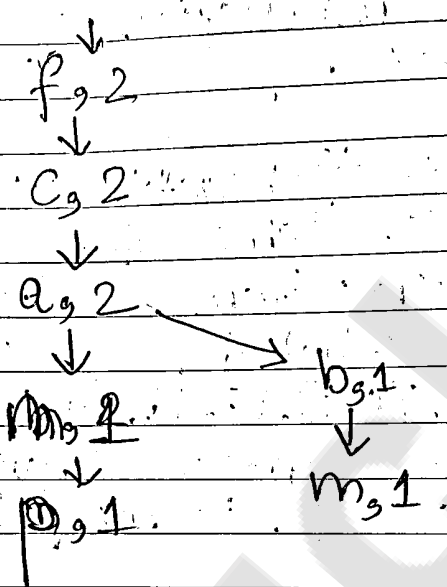


Scan of second transaction  
 $(a, b, c, f, d, m, o)$

$\{(b,1), (c,1), (a,1), (b,1), (m,1)\}$



ROOT

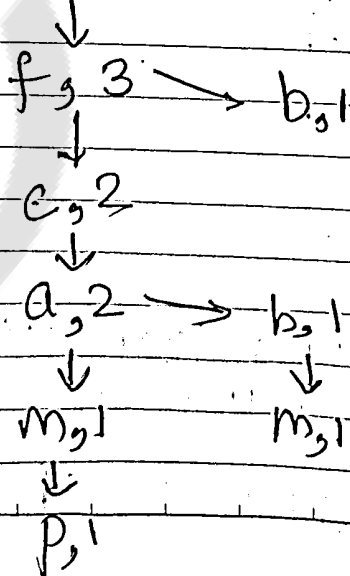


Scan third transaction

$(b, f, h, j, 0)$

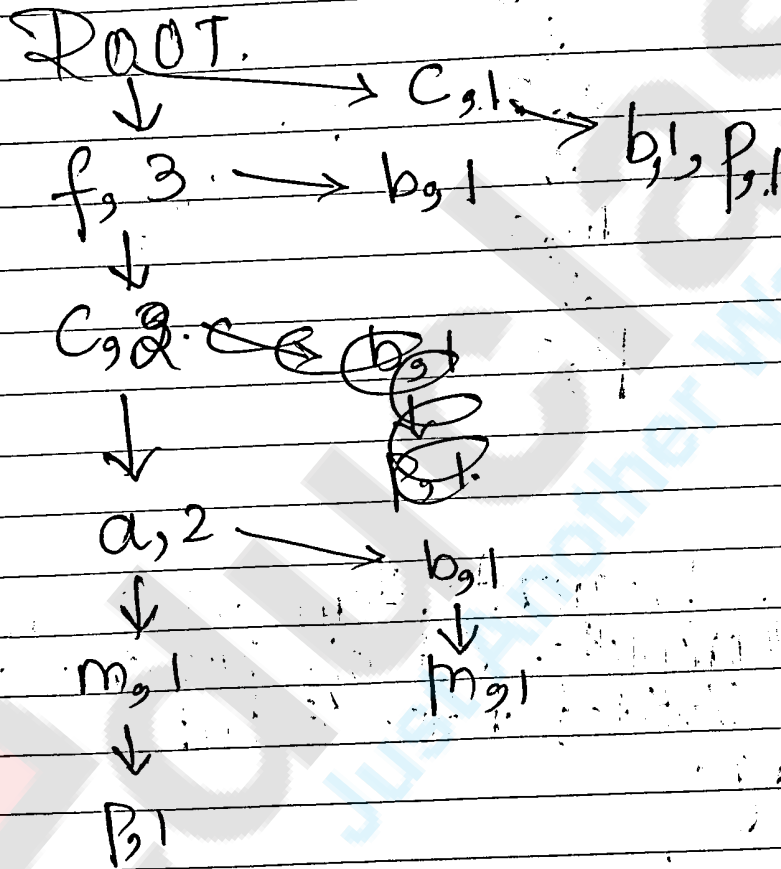
$\{(f_{3,1}), (b_{3,1})\}$

ROOT



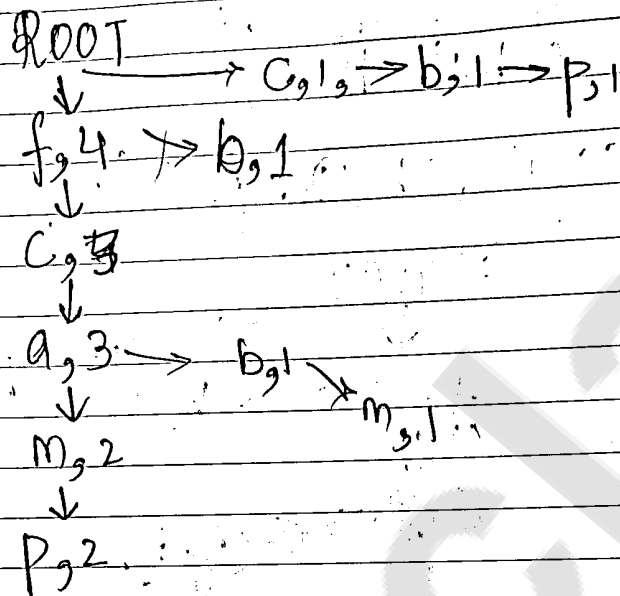
Scan 4<sup>th</sup> transaction  
 $(b, C, R, S, P)$

$\{(C, 1), (b, 1), (P, 1)\}$



Scan 5<sup>th</sup> transaction  
 $(a, f, C, e, l, p, m, n)$

$\{(f, 1), (C, 1), (a, 1), (m, 1), (P, 1)\}$



Finding frequent item set from FP tree.  
The complete set of frequent item set can be divided into 6 subsets without overlapping.

- ① frequent item sets having item p.
- ② frequent item sets having item m but not p.
- ③ frequent item sets having b but not m & p.
- ④ f<sup>a</sup>'s having a without b, m, p.

⑤ freqs having c without a, b, m & p.

⑥ freqs having f.

① freq item sets having p

$\{(f, 4), (c, 3), (a, 3), (m, 2), (p, 2)\}$  and

$\{(c, 1), (b, 1), (p, 1)\}$ .

two parts will be selected from the FP tree i.e.

Samples with a freq item p are

$\{(f, 2), (c, 2), (a, 2), (m, 2), (p, 2)\}$  and

$\{(c, 1), (b, 1), (p, 1)\}$

The given threshold value 3 satisfies only the frequent items are

$\{(c, 3), (p, 3)\}$  or  $\{c, p\}$ .

② freq item set having m but not p.

$\{(f, 4), (c, 3), (a, 3), (m, 2)\}$  and

~~$\{(f, 1), (a, 3), (b, 1), (m, 1)\}$~~

Sample with freq item  $m$  are  
 $\{(f,3) (c,3) (a,3) (m,3)\}$  and  
 $\{(a,1) (b,1) (m,1)\}$

The given threshold  $3$  satisfies  
 $\{(a,3) (m,3)\}$  or  $\{a, m\}$

$\{(f,4) (c,3) (a,3) (m,2)\}$  and  
 $\{(f,4) (c,3) (a,3) (b,1) (m,1)\}$

$\{(f,2) (c,2) (a,2) (m,2)\}$  and  
 $\{(f,1) (c,1) (a,1) (b,1) (m,1)\}$

threshold  $\Rightarrow 3$

$\{(f,3) (c,3) (a,3) (m,3)\}$  or  $\{f, c, a, m\}$

③ freq item set having b.

$\{(f, 4), (c, 3), (a, 3), (b, 1)\}$  and

$\{(f, 3), (b, 1)\}$  and

$\{(c, 1), (b, 1)\}$ .

Sample with freq item b.

$\{(f, 1), (b, 1)\}$  and

$\{(f, 1), (b, 1)\}$  and

$\{(c, 1), (b, 1)\}$

threshold  $\Rightarrow 3$ .  $\{b, 3\}$ .

Since only one item is obtained, it is not considered.

④ f is having a.

$\{(f, 4), (c, 3), (a, 3)\}$

$\{(f, 3), (c, 3), (a, 3)\}$

threshold  $\Rightarrow 3$   $\{b, c, a\}$

⑤ freq item set having  $c$ .  
 $\{f, 4\}, \{c, 3\}$  and  $\{c, 1\}$ .

sample of  $c$ .  
 $\{f, 3\}, \{c, 3\}$  and  $\{c, 1\}$

threshold  $\Rightarrow 3$   $\{f, 3\}, \{c, 4\}$  or  $\{c, 3\}$ .

⑥ freq item set of  $f$ .

Non-frequent item sets we get are  
 $\{c, p\}, \{f, c, a, m\}, \{f, c, a\}, \{f, c\}$ .

$\rightarrow \{f, c, a\}$  and  $\{f, c\}$  are subsets of frequent item set  $\{f, c, a, m\}$ .

∴ the final solution of in FP growth model is:  $\{c, p\}$  and  $\{f, c, a, m\}$

20/02. DMB1

What is the difference between Association and classification?  
Explain: Associative classification method with suitable example.

\* Associative Classification Method. (ACM)

\*\*

CMAR: Classification based on Multiple association rule.

It is a classification association rule mining method which is used for associative classification method.

Method: The association classification method consists of two phases.

- ① Rule Generation / Training
- ② Classification / Testing.

① In rule generation method, CMAR computes the complete set of rules in the form of  $R: P \rightarrow C$  such that  $\text{sup}(R)$  and  $\text{conf}(R)$  pass the given thresholds where  $P = \{a_1, a_2, \dots, a_k\}$  is a pattern with a set of attribute value.  
 $C = \text{class level}$ .

For a given support threshold and confidence threshold the ACM finds



the complete set of class association rules, passing the thresholds.

② In classification phase, when a new sample comes the classifier represented by a set of association rules, selects the rule that match with the sample and has the highest confidence and is used to predict the classification of the new sample.

Q Apply ~~association~~ to the following training data set T with support threshold 2 and confidence threshold 70%

Trd.	Item set	Class.
01	$a_1, b_1, c_1, d_1$	A
02	$a_1, b_2, c_1, d_2$	B
03	$a_2, b_3, c_2, d_3$	A
04	$a_1, b_2, c_2, d_3$	C
05	$a_1, b_2, c_1, d_3$	C

Sol<sup>n</sup>:

Scan the training data set and find the set of item set occurring  $\geq 2$  freq.

ItemSet	frequency
a <sub>1</sub>	4 ✓
b <sub>1</sub>	1
c <sub>1</sub>	3 ✓
d <sub>1</sub>	1
b <sub>2</sub>	3 ✓
d <sub>2</sub>	1
a <sub>2</sub>	1
b <sub>3</sub>	1
c <sub>2</sub>	1
d <sub>3</sub>	3 ✓
c <sub>3</sub>	1

The qualified ItemSets are:

$$\{(a_1, 4) (c_1, 3) (b_2, 3) (d_3, 3)\}$$

Sort the items in descending order:

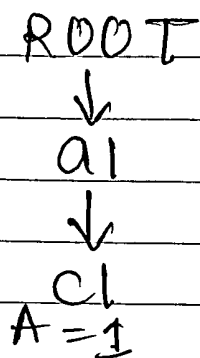
$$F = \{(a_1, 4) (b_2, 3) (c_1, 3) (d_3, 3)\}$$

$$\{a_1, c_1, b_2, d_3\} \Rightarrow \{a_1, b_2, c_1, d_3\}$$

Step 1: Scan the training dataset to construct FP tree.

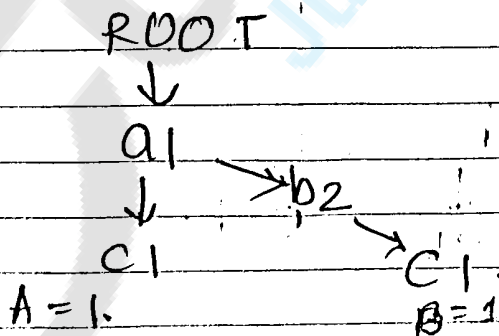
$f(a_1, b_2, c_1, d_3)$

Scan of 1<sup>st</sup> transaction  $(a_1, b_1, c_1, d_1)$ .  
The extracted items are:  $(a_1, c_1)$   
The path is  $a_1 \rightarrow c_1$ ,  
FP tree will be.

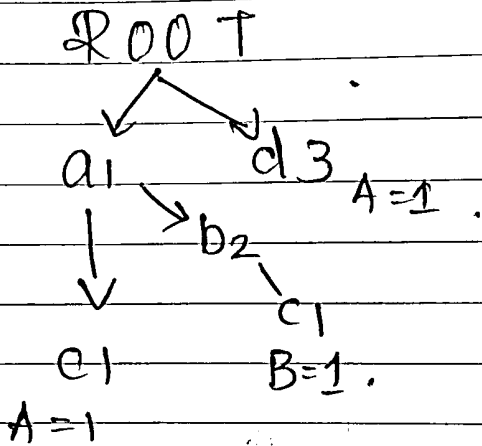


Base of the class label of the sample & the corresponding counter are attached to the last node of the path.  
 $A=1$ .

Scan of 2<sup>nd</sup> transaction  $(a_1, b_2, c_1, c_2)$   
 $(a_1, b_2, c_1)$ ,  $a_1 \rightarrow b_2 \rightarrow c_1$ .  
FP tree will be

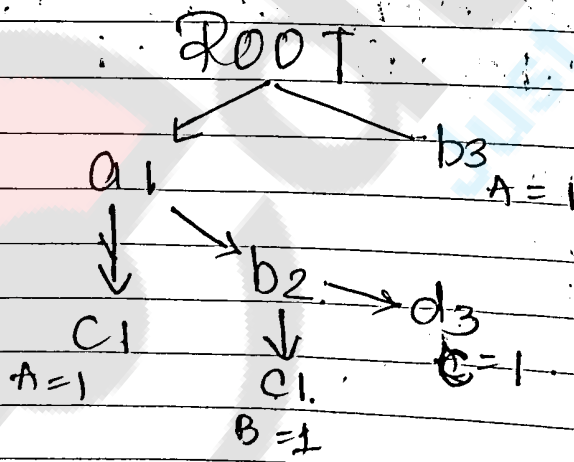


Scan of 3<sup>rd</sup> transaction (a<sub>2</sub>, b<sub>3</sub>, c<sub>2</sub>, d<sub>3</sub>)  
Extracted (d<sub>3</sub>)

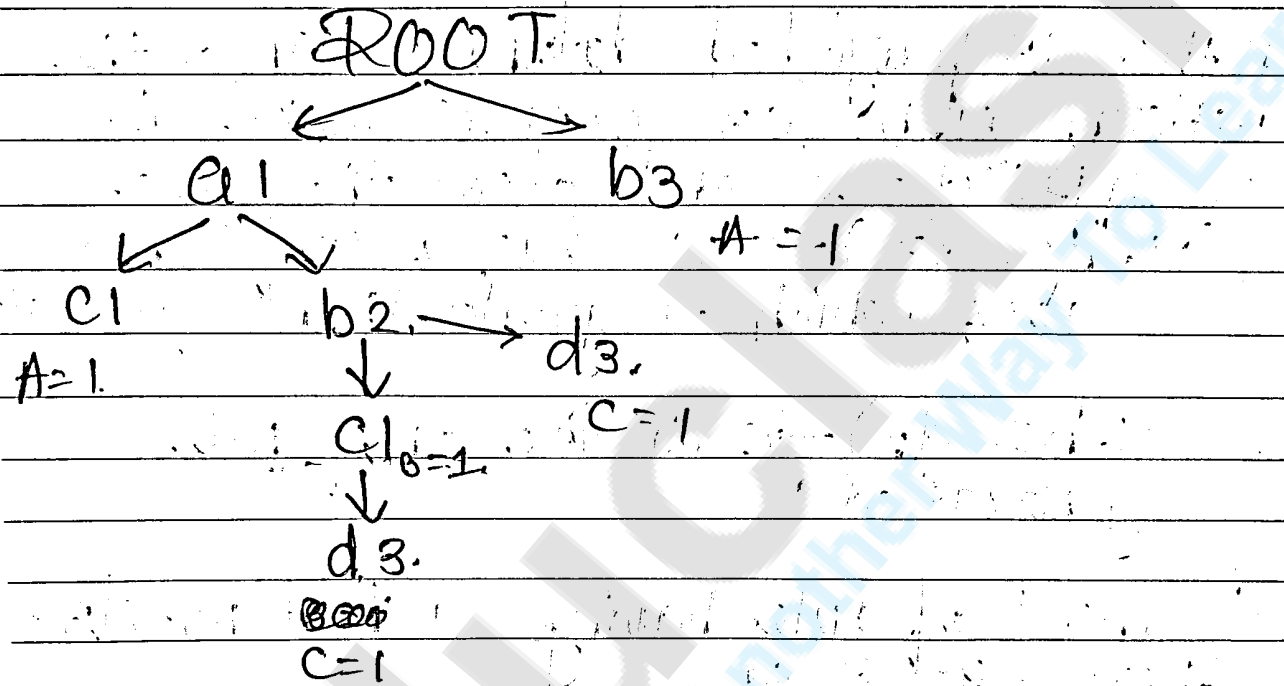


Scan of 4<sup>th</sup> transaction (a<sub>1</sub>, b<sub>2</sub>, c<sub>3</sub>, d<sub>3</sub>)  
Extracted (a<sub>1</sub>, b<sub>2</sub>, d<sub>3</sub>)

a<sub>1</sub> → b<sub>2</sub> → d<sub>3</sub>



Scan of 5<sup>th</sup> transaction.  $(a_1, b_2, c_1, d_3)$   
 Extracted  $(a_1, b_2, c_1, d_3)$   
 $a_1 \rightarrow b_2, \rightarrow c_1 \rightarrow d_3$



Classification Association rule: can be generated by dividing all rules into subsets without overlap.

- ① The rules having  $d_3$  value.
- ② the rules having  $c_1$  but not  $d_3$ .
- ③ the rules having  $b_2$  without  $d_3$  &  $c_1$ .
- ④ the rules having only  $a_1$ .

Subset of rules having  $d_3$  value.  
 there are 3 samples represented in the FP tree with  $d_3$ .

they are  $a_1, b_2, c_1, d_3$  :  $c_1$  | 1. |

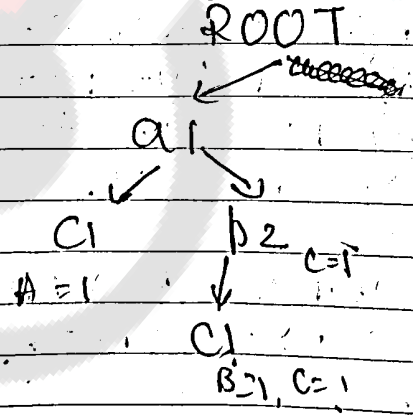
	support	confidence
$a_1 \rightarrow b_2 \rightarrow c_3 : c$	2	$\frac{2}{2} \times 100 = 100\%$
$d_3 : A$	1	

In the db projected database since the pattern  $a_1, b_2, c_3$  occurs twice i.e support = 2 and also the rule  $(a_1, b_2, d_3) \rightarrow c$  has a confidence 100% which is greater than equals 70%

So  $(a_1, b_2, d_3) \rightarrow c$  is the only rule generated

② Subset of rule having  $c_1$  but no  $d_3$ .  
~~there are 2 sa~~

After a search for rules having  $d_3$  value all the nodes of  $d_3$  and their corresponding class labels are merged into their parent nodes at the FP tree.

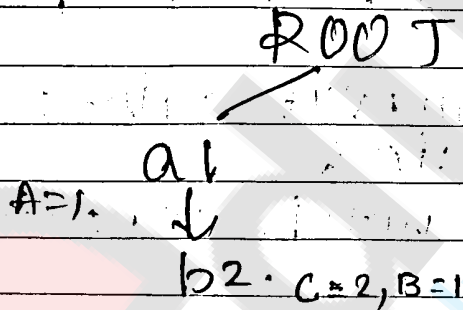


There are 8 samples represented in FP tree with  $c_1$

	Support	Confidence
$a_1 c_1 : A$	1	1
$a_1 b_2 c_1 : B$	1	1
$a_1 b_2 c_1 : C$	1	1

all the rules are with  $\leq 2$  value, so no rule will be generated.

③ Subset of rule having  $b_2$  but not  $a_3, c_1$ .



There are 2 samples with  $b_1$

$a_1, b_2 : G$	2	$\frac{2}{3} \times 100 = 66.66\%$
$a_1, b_2 : B$	1	

rule  $(a_1, b_2) \rightarrow C$  has support value 2 but confidence  $< 70\%$

So, no rule will be generated

① Subset with a  $\bar{q}$ .

Root



$a_1$

$A=1 \ B=1 \ C=2.$

3 sample

$a_1 \rightarrow A$

1

$a_1 \rightarrow B$

1

$a_1 \rightarrow C$

2

$\frac{2}{4} \times 100 = 50\%$

$(a_1) \rightarrow C$  has support value 2

but confidence  $< 70\%$

So this rule will not be considered.

$\therefore$  the only association rule generated from the trending process with the above given db is  $(a_1, b_2, d_3) \rightarrow C$  with support value 2 & confidence 100%.

cont. back



## DMB1

phase 2:

- Classification or testing:  
for a new sample CMAR collects the subset of rules matching the samples from the total set of rules.
- If all the rules have same class CMAR simply assigns that level to the new sample.
- If the rules are not consistent in the class level CMAR divides the rules into groups according to the class level & find the strongest root.
- CMAR uses the strongest rule in the group as its representative i.e. the rule with highest  $\gamma^*$ .