

Unit - 1

- • Overview
- • Motivation [Data Mining]
  - Data Mining - Definition and functionalities
  - Data Warehousing
    - overview
    - definition
    - delivery process
    - difference between Database System & Data Warehouses
  - Multidimensional data Model.
  - Data Cubes
  - Snow flakes
  - Fact Constellation
  - Concept Hierarchy
  - Process Architecture
  - Three Tier architecture
  - Data Marting
    - ROLAP
    - MOLAP
    - NOLAP

Unit - 2

- Data Preprocessing
- Data Cleaning : Missing Values : noisy data [binning]

- clustering
- Regression
- Inconsistent data.
- Data Integration and transformation.

→ Data Reduction:

- data cube aggregation
- dimensionality reduction
- data compression
- numerosity reduction
- CLUSTERING
- Discretization
- Concept Hierarchy Generation.

### Unit-3

Concept Description

- Definition
- Data Generalisation
- Analytical characterisation
- Analysis of attribute relevance.

→ Mining class comparison:

→ Statistical measure in large database.

→ Measuring Central Tendency.

→ Measuring Dispersion of data.

→ Apriori Algorithm.

→ Mining Multilevel Association rule from transactional database.

Unit-4:Classification

- What is classification.
- Issues regarding classification.
- Decision tree.
- Bayesian Classification.
- Classification by back propagation

Unit-5:Cluster Analysis

- Data types in cluster analysis.
- Partitioning method.
- Hierarchical Clustering  
CURE and Chameleon
- Density based methods.  
→ DBSCAN → OPTICS.
- Grid based methods  
→ STING → CLIQUE
- Outlier Analysis

→ Difference b/w database & data warehouse

↓  
collection of  
records

Definition: Data warehousing is a relational database management system (RDBMS) designed specially, to meet the needs of transaction processing system.

RDBMS: It is a database, that stores information oriented, satisfy decision making request.

Father of Data Warehousing: Bill Inmon

He provides a definition as:

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

Features of Data warehouse:

- It is separate from operational database;
- It integrates data from heterogeneous system.
- It does not require data to be highly accurate.
- Queries are generally complex.
- Supports online analytical processing system [OLAP]

: Influence decision making in favor of the Enterprise.

### Goals of Data Warehousing

- Data warehousing maintains organisations' historical information.
- Data warehousing is said to be foundation for decision making.
- Helps in reporting and as well as in analysis.

### Need of Data Warehouse

#### (1) Business User

It requires data warehouse to view summarised data from the past. [for simpler representation].

#### (2) Stored Historical Data

Data warehouse is required to store the time variable data from the past.

#### (3) Make Strategic Decision

Some strategies may be depending upon the data in the data warehouse.

## ② High Response Time

Data warehouse has to be ready for rarely unexpected loads and types of queries, which demand a high degree of flexibility and quick response time.

## ③ For Data Consistency and Quality

Bringing the data from different sources, at a common place, user can effectively undertake to bring the uniformity and consistency in data.

## Difference between operational Database System and Data Warehouse

Features	Operational DB	Data Warehouse
Users	(OLTP)	(OLAP)
Workload	→ thousands → present transactions	→ hundreds → specific analysis queries
Access	→ to hundred of records	→ to millions of records
Goal	depends on decision making support application	→

Data Integration → Application based → subject based

Quality → In terms of integrity → In terms of consistency

# OLTP: Online Transaction Processing System.

Charlie ①

Date \_\_\_\_\_

Page No. \_\_\_\_\_

- Time coverage → Current data only → current and historical data
- Updates → continuous → periodical
- Model → Normalised → denormalised
- Optimization → for OLTP access → for OLTP query to db part most of db.

Operational Activities: It may be categorised in OLTP system.

Data warehouse: It may be categorised in OLAP system.

\* The data warehouse and OLTP database are both relational database, but the objective of both these database are different.

## DATA WAREHOUSE Db

② Designed for analysis of business measure by categories and attributes

## OLTP Database

Designed for real time business operations.

③ Optimized for bulk loads and large, complex, unpredictable queries

Optimized for a common set of transactions usually adding or retrieving a

that access many rows per table

single row at a time per table.

③ loaded with consistent, valid data, requires no real time validation.

Optimized for validation of incoming data during transactions

④ supports few concurrent users relative to OLTP

Supports thousands of concurrent users.

Advantages of OLAP

disadvantages of OLAP

## Differences between OLTP and OLAP

Features	OLTP	OLAP
1) Source of Data	Operational Data	Consolidation Data
2) Purpose of Data	To control and run fundamental business class.	To help with planning, problem solving and decision support.
3) What the data is?	Receives a snapshot of ongoing business process	Multidimensional views of various kind of business activities
4) Stored values.	Strictly with the coded data	Stores descriptive data.
5) Processing Speed	Fast	Depends on amount of data.
6) Queries	Standardise and simple query.	Complex queries.
7) Space Requirement	Small	Large
8) Database Design	Highly normalised with many tables	De-normalised with few tables.
9) Normalisation	Fully normalised	Partial normalised
10) Modeling	Very entity relation model.	Dimensional model snowflake

ii) User

Clock and IT professionals.

Knowledge worker

iii) Database Design

Application oriented.

Subject Oriented.

OLTP : Online Transaction Processing System.

The major task of Online Operational database system, is to perform online transactions and quick processing.

The cover day to day operations such as  
→ purchasing, inventory, registration and accounting.

OLAP : Online Analytical Processing System.

Data warehouse, on the other hand serve user or knowledge workers in the role of data analysis and decision making.

Such systems can organise and present data in various formats in order to accumulate the advise needed by the different users.

- Clustering
- Classification.

Date \_\_\_\_\_  
Page No. \_\_\_\_\_

11  
Charlie

## Data Mining

It is a process of extracting hidden knowledge from large volume of raw data, and using it to make crucial business decision.

### Functions of Data Mining

#### 1) Classification

→ Infers the defining characteristics of a certain group.

#### 2) Clustering

→ Identifies groups of items that share a particular characteristics.

#### 3) Association

→ market strategy

→ Identifies relationship between events that occur at one time,

Such as: shopping basket.

Bread + Jam

#### 4) Segmentation

## ④ Scanning

Similar to association rule, except that the relationship ~~exists~~ exists over a person.

### Some of the benefits

#### (1) Fraud Management

Telecommunication, financial, insurance industries.

#### (2) Market Analysis

Customer, Competition

#### (3) Entertainment

Digital convergence

#### (4) Diagnosis and Monitoring

Medical

Q. Why do we use Data Mining?

(1) Human analysis, scalar inadequate:

- i. Volume and dimensionality of data.
- ii. High data growth rate.

(2) Availability of:

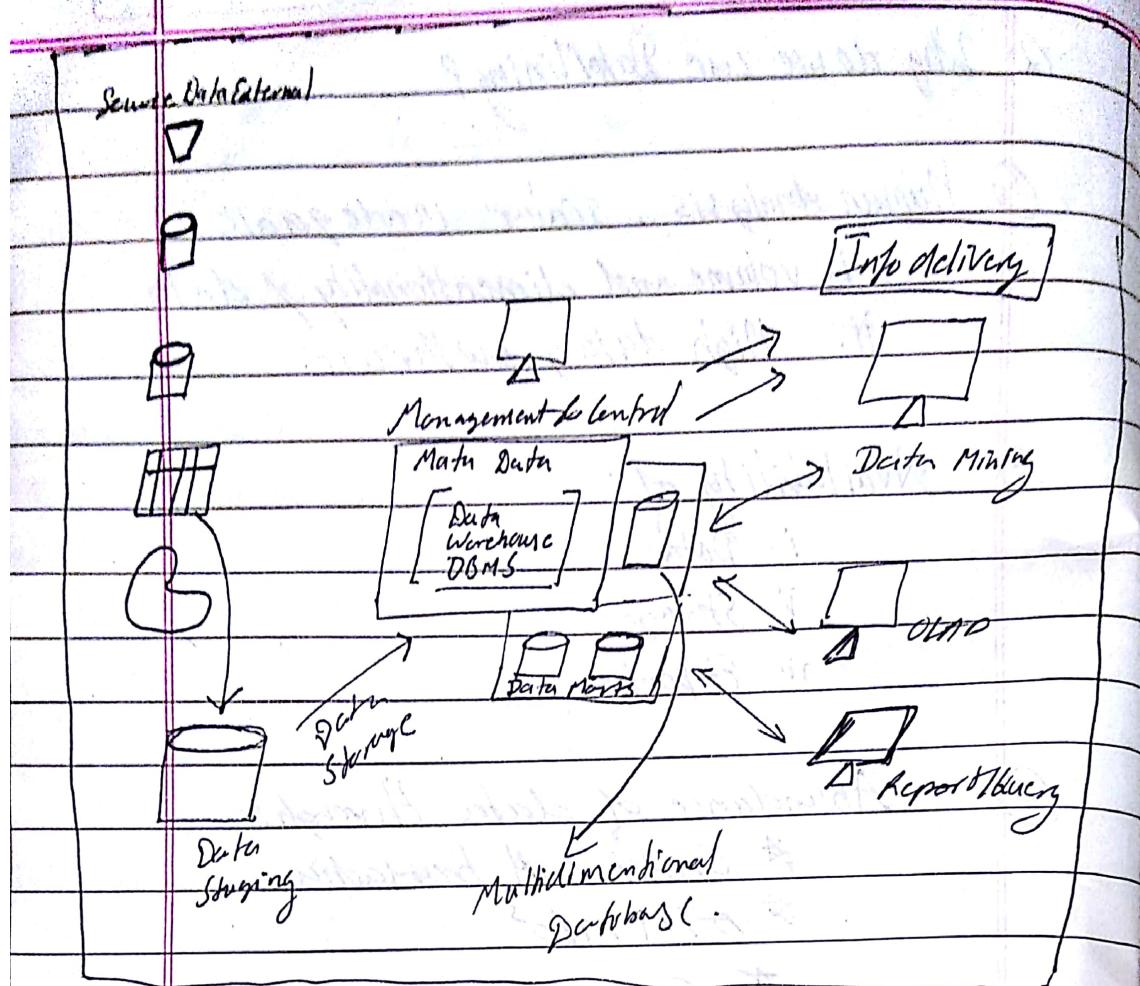
- i. Data
- ii. Storage
- iii. Expertise

(3) Abundance of data through:

- \* Credit card transactions
- \* ATM m/s
- \* Camera
- \* Direct Mail Response
- \* Call center records
- \* Sensor Networks.

(4) Forecasting:

- \* Estimates future values, based on patterns within largest set of data.



Data Warehouse Building Blocks or Components.

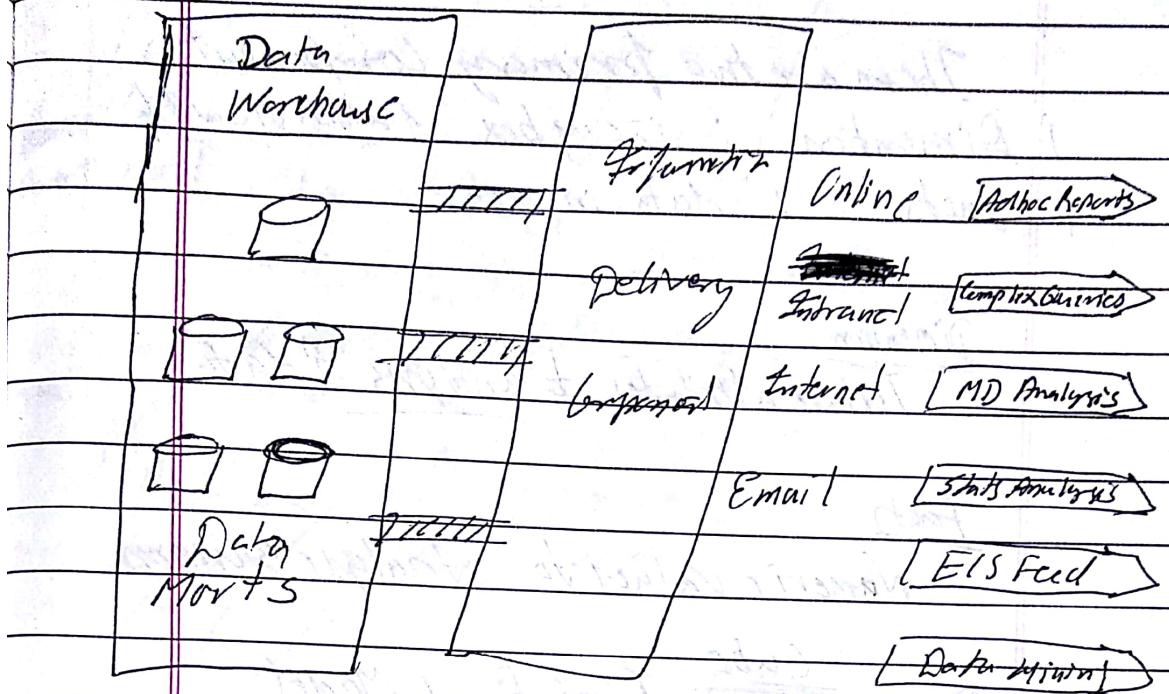
Data warehouse : collecting records

Data marts : Catalogue.

15  
Charlie

Date \_\_\_\_\_  
Page No. \_\_\_\_\_

## Information Delivery Component



## Multi Dimensional Data Model

(MDDM)

The dimensional model was developed for implementing data warehouse and data Marts.

It is an integral part of online analytical processing (OLAP)

This is designed to solve complex queries in real time

## Component of Multi-Dimensional Data Model

There are two primary component.

- i. Dimension : size of box of which view?
- ii. Facts : data in cells

Dimension

Texture attributes to analysis on data.

Facts

Numeric values to analyse Business

DATA ~~FACT~~ Cube Dimensional Model

When data is grouped combined together in multidimensional matrices, then it is called data cube.

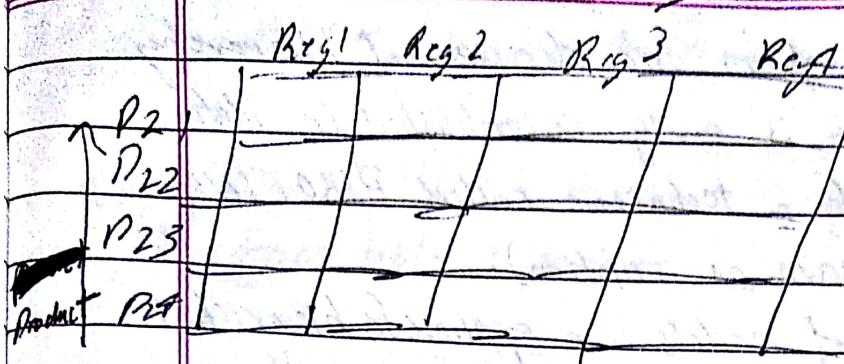
In Two Dimension: There's row and column

or

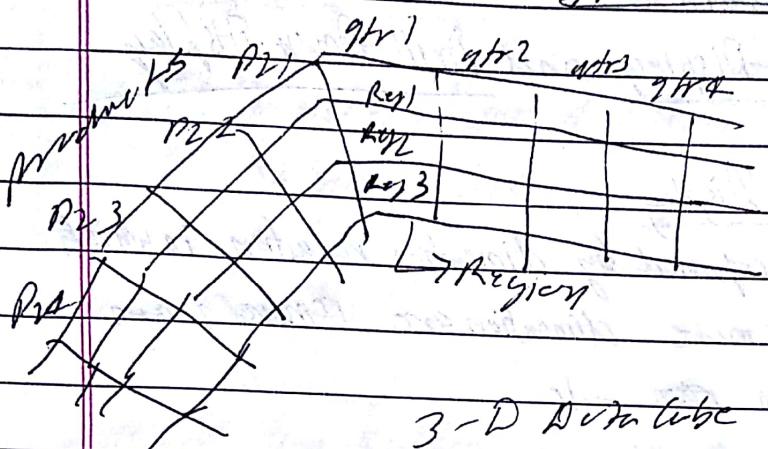
Product and ~~Region~~ fiscal quarters

In 3-D: one region, products and Fiscal Quarter

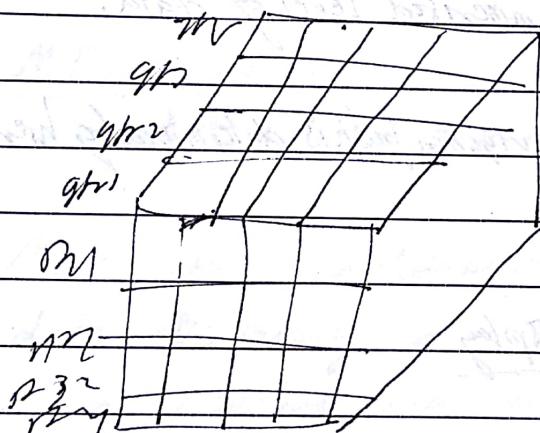
## Regions



Fiscal quarter



3-D Data cube



Changing from one dimensional hierarchy to another is easily accomplished in data cube by a technique called PIKOTING (also known as rotation).

These type of models are applied to hierarchical view, such as:

Roll-Up-Display and Drill-Down-Display

Roll Up Display

It is performed by dimension reduction in which one or more dimension are removed from dimension ~~from~~ role.

With the role of capability user can zoom out to see a summarised level of data.

The ~~the~~ navigation path is determined by hierarchy in dimension.

Drill-Down-Display

- It is reverse of Roll Up Display.
- It ~~navigates~~ moves from less detailed data to more detailed data.
- It can also be performed by adding new dimension to a cube.

MMDM involve two types of tables:

- i. Dimension Table
- ii. Fact table.

### Dimension Table:

- It consists of tuples of attributes of dimension.
- It is a simple primary key.

### Fact Table

- It has tuples, one per fact recorded fact.
- It is a compound primary key.

### META DATA

- Information about data.

Metadata in a data warehouse is similar to data dictionary or the data catalogue in a database management system.

### Types of Meta data.

#### (1) Operational Meta Data

#### (2) Extraction and Transformation Meta Data.

#### (3) End User Meta Data.

## Operational

It contains all kind of information, about the operational data sources.  
(Split, Add etc)

## Extraction and Transformation of Meta Data

It contains data about the extraction of data from the source system, namely:

Extraction Frequency

Extraction Method

It also contains information about all the data transformation that take place in the data staging area.

## End User Meta Data

It is the navigational map of the Data Warehouse.

It enables End User to find information from the data warehouse.

## Goal of Multidimensional Data Model

- It supports analysis in a simple and faster way for executives, managers and business professionals.
- Application:
  - Used in building online analytical processing engines.

## Data Cube

When data is grouped, or combined together in Multidimensional Matrices, then it is called Data Cube.

If it is constructed from a subset of attributes in the database.

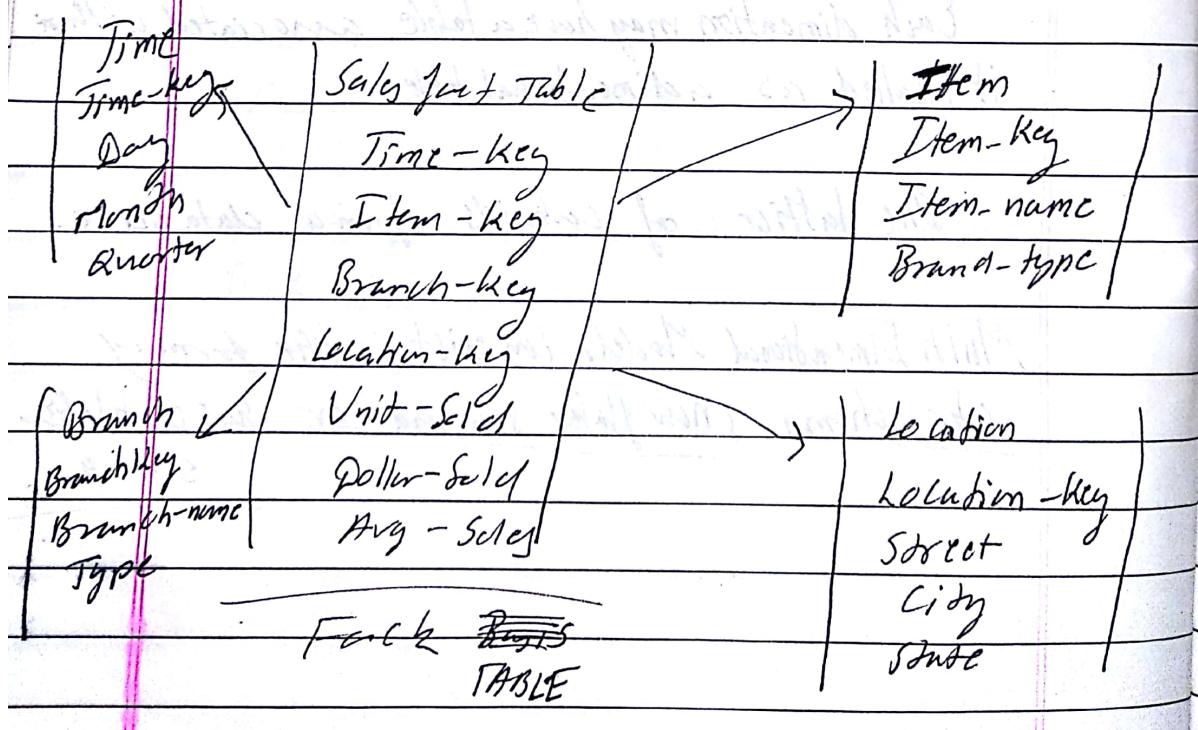
Each dimension may have a table associated with it, called as a dimensional table.

The lattice of cuboids form a data cube.

Multi Dimensional Model, can exist in the form of Star schema, Snowflake schema or fact constellation schema.

## STAR Schema Model

- It is also known as STAR join Schema
- It is the simplest style of Data warehouse schema.
- It is called 'star schema' because ~~its~~ Entity relationship diagram of this schema resembles a STAR with points radiating from the central table.
- A star query is a join between a fact table and a number of dimensional table.
- Each dimension table is joined to the fact table using primary key to foreign key join, but dimension table are not joined to each other.



## Characteristics

- It creates a denormalised database, that can quickly provide query responses.
- It reduces the complexity of Relational Data for both developers and End User.

## Advantage

- It has simple structure.
- It has smaller no. of tables and clear join paths.
- Easy for end users and applications to understand and interact.

## Snow Flake Schema

- It is an extension of Star Schema
- It is different from Star Schema, in which the dimension tables from a star schema are ~~are~~ organised into a hierarchy by normalising them.
- It is represented by Centralised fact table which are connected to multiple dimension
- The snow flaking effecting <sup>only</sup> the dimension table not the fact table

### \* Definition

Snow flaking is a method of ~~not~~ normalising the dimension table in a star schema.

It improves the performance of certain queries.

### Advantages

- No ~~any~~ redundancy
- Greater flexibility

### Disadvantages

- More complex queries
- Difficult to understand.
- More tables, more query execution time.

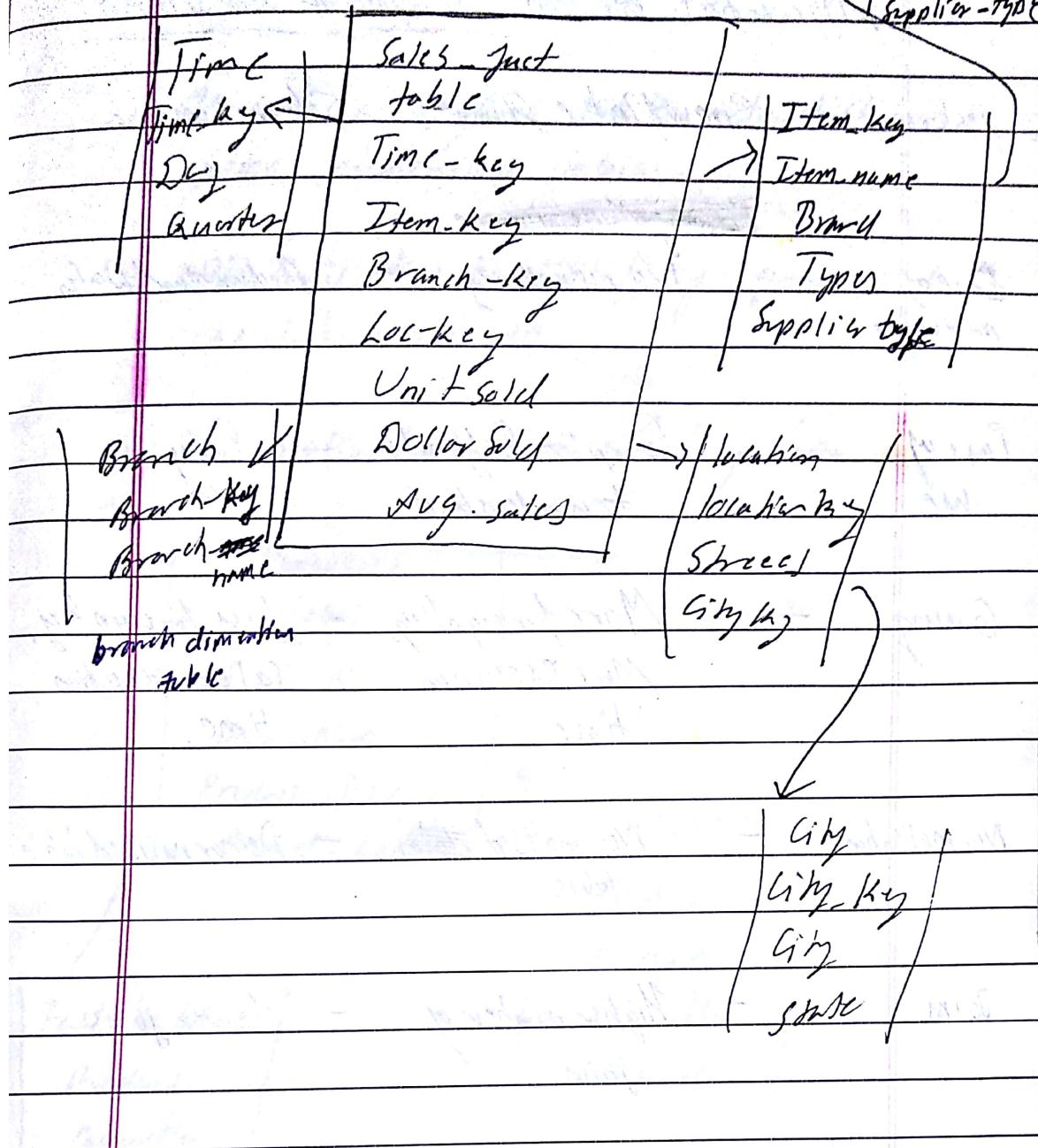
28

Charlie

Date

Page No.

Supplier  
Supplier-key  
Supplier-type



## Difference between Snowflake & Star Schema

Characteristics - Snowflake Schema - Star Schema

~~Easy Maintenance~~  
Ease of maintenance - No redundancy - Redundant Data.

Complexity - Complex, Difficult to understand - Easy

Query - More foreign key, - less foreign key.  
More execution time - So less execution time.

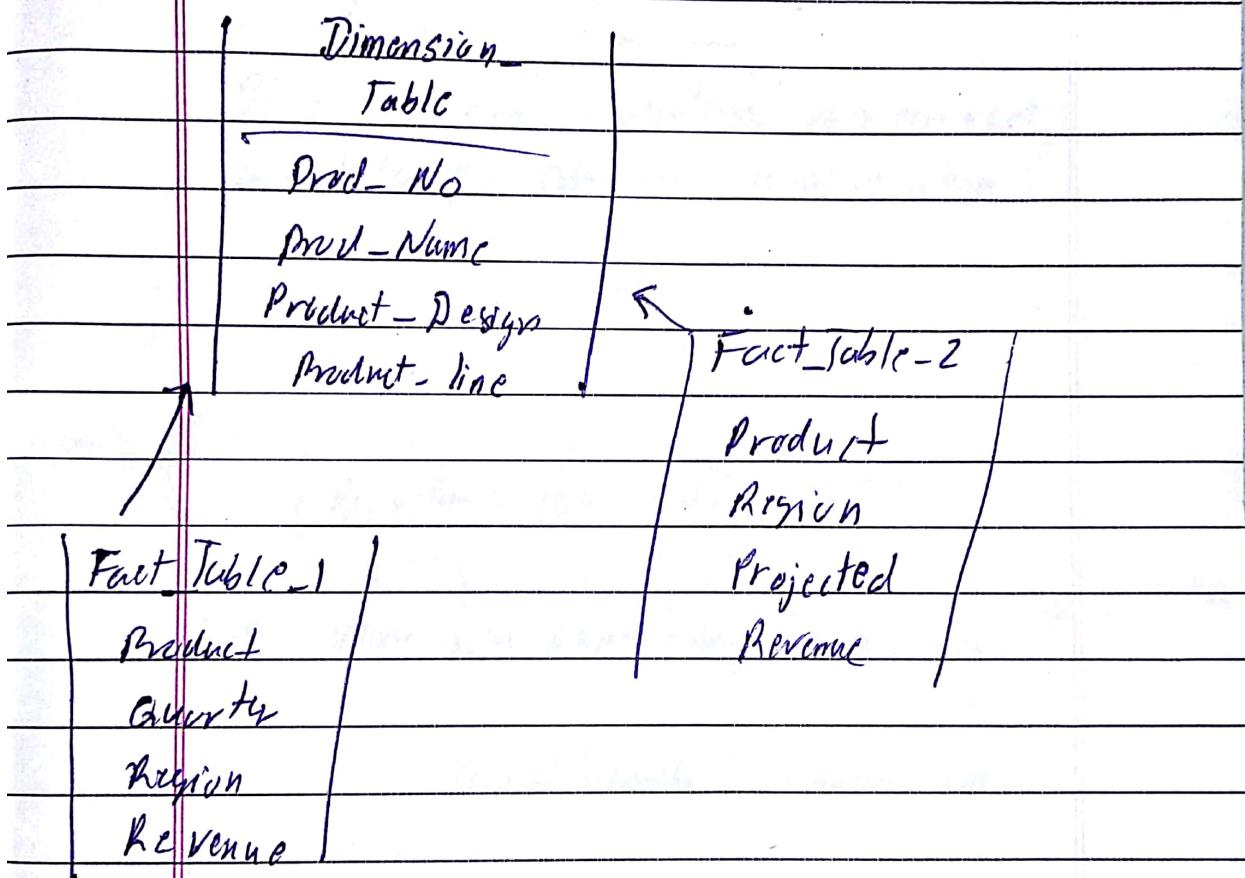
Normalisation - Normalised ~~table~~ - De-normalisable table

Join - Higher number of joins - Fewer joins

Data warehouse systems - Better for small data warehouse, data mart - Works best in any data warehouse, data mart

## Fact Constellation

- i. It is a set of fact tables, that share some dimension tables.
- ii. It limits the possibility of querying for the data warehouse.
- iii. It is also known as Galaxy Schema.



## Assignment - 1

- ① → Discuss the delivery process in Data warehouse
- ② → Discuss ROLAP, MOLAP and HOLAP
- ③ → Difference between data Warehouse and Data Mart

→ Concept Hierarchy

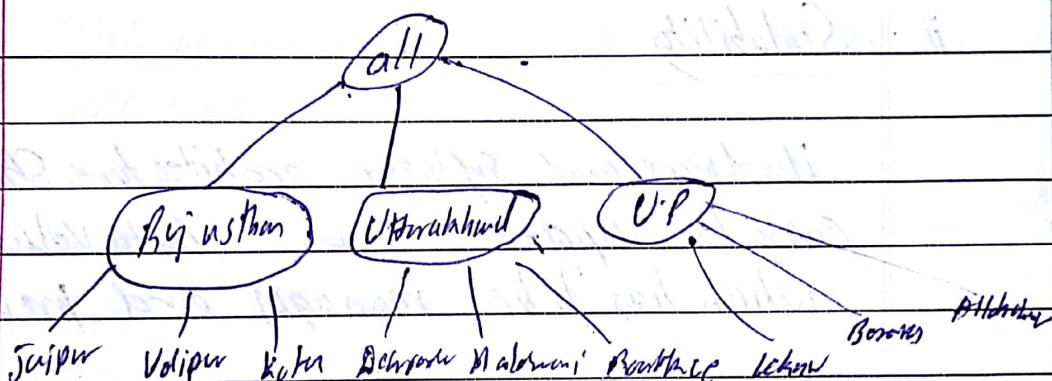
→ Process Arch.

→ 3-Tier Arch

### Concept Hierarchy

It is defined as a sequence of mapping from a set of low level concepts to higher level concepts to a more generalised concept.

It is an order relation, between a set of attributes of a concept or dimension.



Concept hierarchy for dimension location

Concept hierarchy may also be designed for grouping values for a given dimension or attribute, resulting in a set of grouping hierarchy.

It can be provided manually by system users, domain expert or knowledge engineer.

## Data Warehouse Architecture

There are some properties for a data warehouse system:

### i. Separation

of OLAP, OLTP

Insertion  
Deletion  
Update

in shared  
in database

Analytical and transactional processing  
Should be kept apart as much as  
possible.

### ii. Scalability

Hardware and software architecture should be  
easy to upgrade as data volume  
which has to be managed and provided.

### iii. Extensibility

The architecture should be able to host  
new application and technology  
without redesigning the whole  
system.

#### iv. Security

Monitoring Access is essential, because of the Strategic data stored in data warehouse.

#### v. Administrability

Data warehouse management, should not be overall difficult.

### 3 Tier Data Warehouse System

Data warehouse adopts 3-tier architecture:  
There are :

- ① Bottom Tier [Data Warehouse Server]
- ② Middle Tier [OLAP Server]
- ③ Top Tier [Front End Tools]

#### Bottom Tier

- It acts as a relational database system,
- Data from operational databases and external sources are extracted using application program interface known as Gateways. [Such as : customer Profiler, information provided by external consultant]
- A gateway is supported by the underlying

DBMS and allow client program to generate SQL code to be executed at the server.

### Middle Tier

It is an OLAP server that is typically implemented by using either of the following:

- ~~An external relational OLAP [ROLAP].~~
  - It is an ~~extended~~ external relational database management system, that maps operations on multidimensional data to standard relational operation.
- ~~Multidimensional OLAP. [MOLAP].~~
  - It is a model that is a special purpose server, that directly implements multidimensional data and operations.

### Top Tier

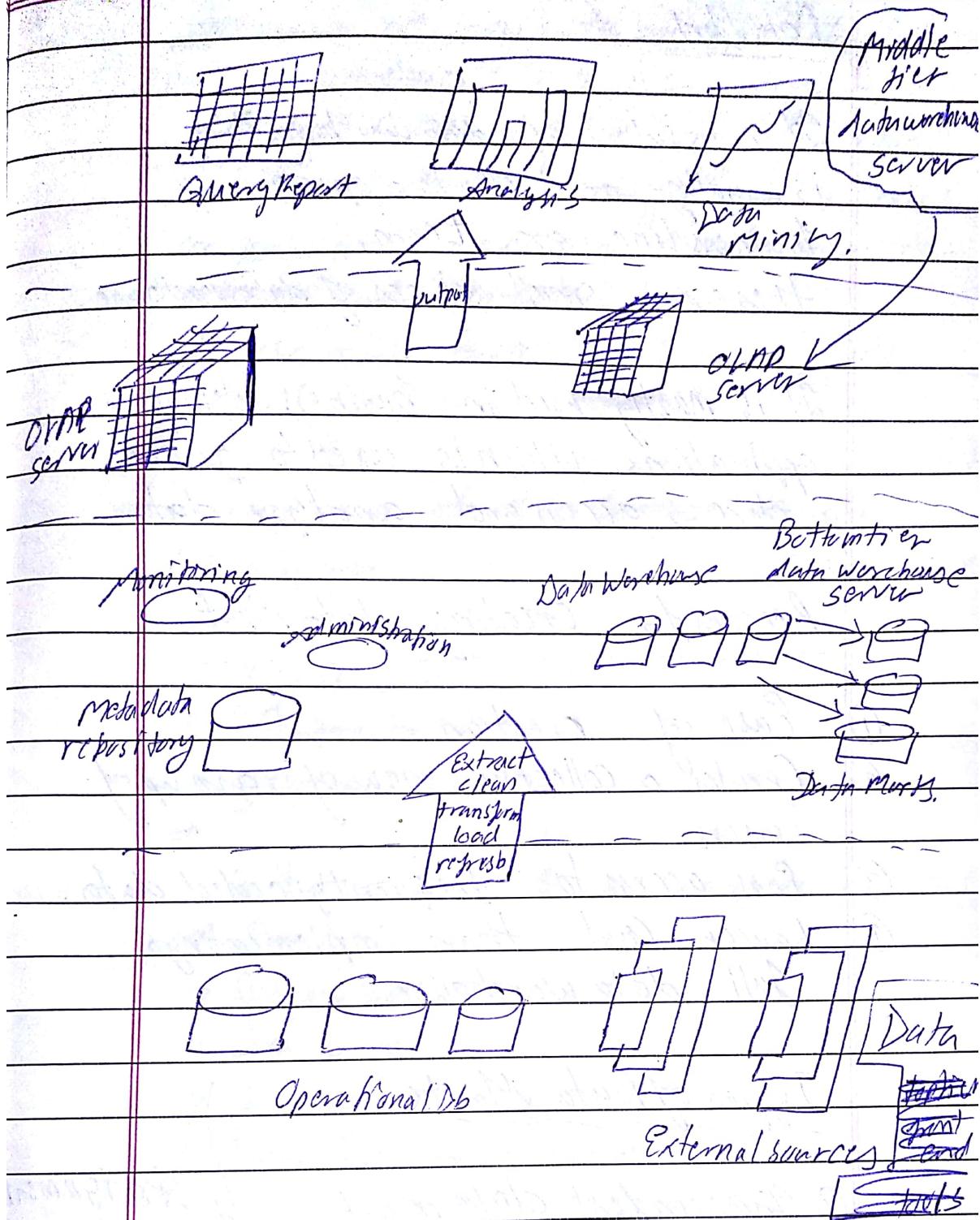
- It is a client which contains query and reporting tools, analysis tool or data mining tools.

eg: Trend Analysis, prediction and so on.

TOP Tier  
Front end  
URLS

Charlie  
33

Date \_\_\_\_\_  
Page No. \_\_\_\_\_



3 Tier Data Warehouse

Architecture

## Data Mart

It is a subset of data warehouse that is usually oriented to a specific business line or team.  
These are small slices of data warehouse.

It is mostly used in business intelligence applications, which is used to gather, store, access and analyse data.

### Reasons for creating data Mart

- ① Ease of creation.
- ② Creates a collective view of group of users.
- ③ Easy access to frequently needed data.
- ④ Lower cost than implementing a full data warehouse.

### Types of Data Marts

- ① Independent data mart
- ② Dependent data mart.

→ Architecture  
→ Advantages & Disadvantages  
→ Differentiation.

The data we wish to analyse by data mining technique, are in incomplete form, (lacking attribute values or only certain attribute of interest), noisy (containing errors or outliers value which deviate from the expected and inconsistent data).

Discrepancies error in departmental code, used to category data items

Why we need ~~to process~~ the data!

- Redundancy
- Aggregation
- Removal of Noisy Data
- Inconsistent Data
- Generalisation

② Why we need ~~the~~ transformation of data!

- Durability
- Data analysis
- Data scaling
- Consistency
- Turn data into a simple form

Different types of Data Processing techniques.

- Data Cleaning
- Data Integration
- Data Transformation
- Data Reduction

## Data Cleaning

It can be applied to remove noise and correct inconsistent data.

## Data Integration

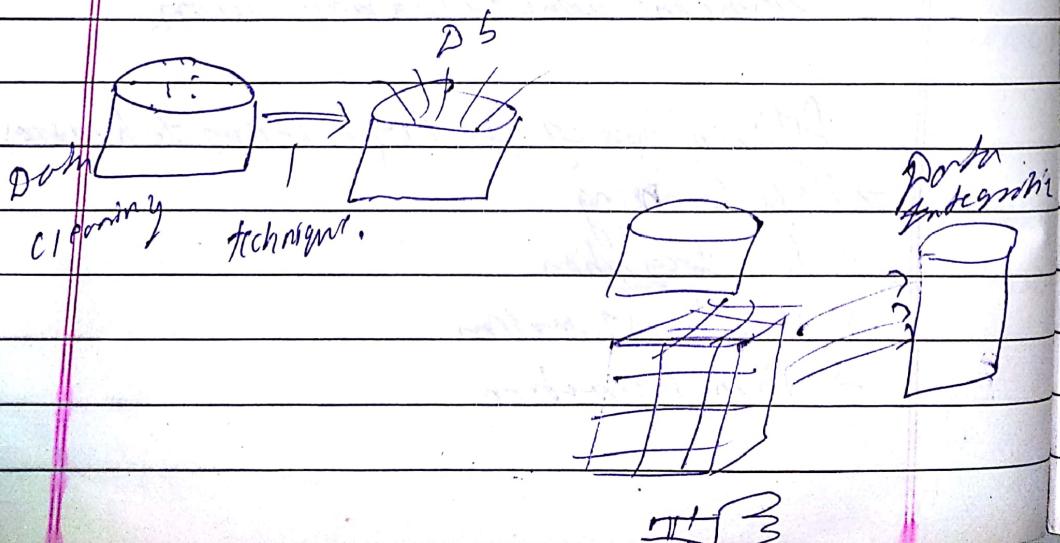
It merges data from multiple resources, into coherent data.

## Data transformation

In this technique normalization can be applied to improve the accuracy, efficiency of Data Mining Algo.

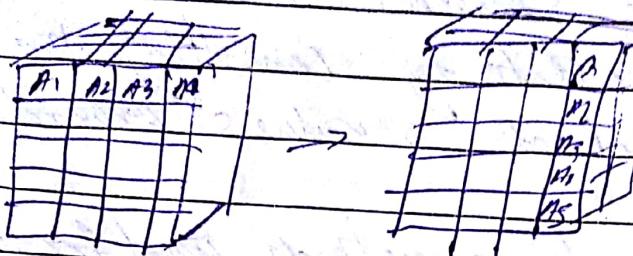
## Data Reduction

It reduces the data size, by aggregating, eliminating redundant data or clustering.



Front view  $\rightarrow$  top view  
 $\rightarrow$  similar

## Data Reduction



## Noisy Data

It is a random error, or a variance in a measured variable, due to randomness, it is very difficult to follow the strategy for noise removal from the data.

Common methods for removing noise from data, are as follows:

- (1) Binning Method.
- (2) Regression
- (3) Clustering
- (4) Combine Computer & Human Inspection.

## Binning method

It smooths the data or sorts the data by consulting the neighbourhood or values around it.

There are three methods used for binning.

- i. Partition into equal frequency bins.
- ii. Smoothing by bin means.
- iii. Smoothing by bin boundaries.

### Numerical

→ Let there be the data price in \$,  
4, 8, 15, 21, 21, 24, 25, 28, 34

Solve the question by binning method:

- i. Partition into equal frequency bins.
- ii. Smoothing by bin means.
- iii. Smoothing by bin boundaries.

(1) Bin 1 (4, 8, 15)

Bin 2 (21, 21, 24)

Bin 3 (25, 28, 34)

(a) Bin 1 (9, 9, 9) Mean of bin 1 way

Bin 2 (22, 22, 22)

Bin 3 (29, 29, 29)

③ Replace by minimum.

Bin 1 (9, 9, 15)

Bin 2 (21, 21, 24)

Bin 3 (25, 25, 34)

## REGRESSION

- Regression data can be smoothed by fitting the data to a function
- Linear regression involves :  
finding the best line to fit the values / two variables so that one variable can be used for predicting the other.
- Multiple linear regression.  
it is an extension of linear regression where more than two variables are involved, and the data are fit to a multidimensional surface.

regression is used for -

- Classification, prediction of data.

Date	Charlie
Page No.	.

## CLUSTERING

In clustering, grouping of data into different groups are performed so that data in each group shares similar trends and patterns.

## Data Integration and transformation

Data integration means, merging of data from multiple data source stores.

This data may also be needed to transform into forms  $\Theta$  appropriate for mining.

The source may include multiple databases.  
The data integration problem:

i.e Redundancy, can be removed by  
Co-relation.

## Numerical

Find the Co-relation!

No. of study hours

(x)

No. of sleeping hours

(y)

2

10

4

9

6

8

8

7

10

6

3  
3040

$$Exy = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \cdot \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

$$= \frac{-20}{20} = 1$$

$\bar{x} = 6$	$x - \bar{x}$	$(x - \bar{x})^2$	
-4	16		
-2	4		$6(x - \bar{x})(y - \bar{y})$
0	0		-8
2	4		-2
4	16		0
	<u>40</u>		-2
			-8

$\bar{y} = 8$	$y - \bar{y}$	$(y - \bar{y})^2$	
2	4		<u><math>\sum -20</math></u>
1	1		
0	0		
-1	1		
-2	4		
		<u>10</u>	

### Note

If the resulting value = 0, then we say that  $X$  and  $Y$  are independent.

If the resulting value is less than 0,

then we say that  $X$  and  $Y$  are negatively correlated, where the values of one attribute, increases as the values of other attribute decreases.

### Data Transformation

In data transformation, the data are transformed into appropriate forms for mining.

The data transformation can involve, the following:

- i. Smoothing
- ii. Aggregation
- iii. Generalisation

## Smoothing

Its function is to remove noise from the data.

These techniques involve:

- binning
- regression & clustering

## Aggregation

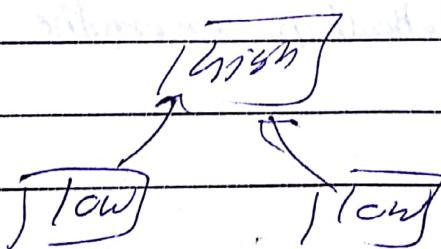
In this, the summary or aggregation operation are applied for the data,

e.g.: the daily sales data may be aggregated so as to compute monthly and annual total amount.

## Generalisation

Generalisation of data, where low level or primitive data are replaced by high level concept, through process of conceptual hierarchies, e.g.: 'Categorical attributes'

like street can be categorical, on higher level concept, like City or Country.



## Normalisation

In this the attributes are scaled, so as to fall within a small specified range; there are many methods for normalisation:

- i. Min-Max Normalisation

→ it performs linear transformation on the ~~the~~ original data.

Suppose that:

Min A and max A are the minimum and maximum values of an attribute then, min-max can be computed as:

$$V' = \frac{V - \text{min}A}{\text{new\_max}A - \text{new\_min}A}$$

where  $V$  to  $V'$  is the range and  $\text{new\_max}A$ ,  $\text{new\_min}A$  represent the relationships among the original data values.

Numerical

Suppose that, the minimum and maximum values for the attribute income are £12,000 and £98,000 respectively.

We would like to map income to the range [0, 1.0] by Min-Max normalisation. A value of £73,600 for the income, will be transformed to?

$$v' = \left[ \frac{73600 - 12000}{98000 - 12000} \right] \cdot [1.0 - 0] + 0$$
$$= \frac{61600}{86000} \cdot 1$$

$$= \frac{616}{860} \cdot 1$$
$$= 0.716$$

## Z-score Normalisation

In this the values for an attribute  $A$  are normalised based on the mean and standard deviation of  $A$ .

A value ' $v$ ' of  $A$  is normalised to ' $V'$ ' &

is given by  $\frac{v - \bar{A}}{\sigma_A}$ , where  $\bar{A}$  =

and  $\sigma_A$  are mean and standard deviation.

This method of normalisation is useful when the actual minima and maximum attribute  $A$  are known.

### Numerical

a. Suppose that the mean and the standard deviation of the values for the attribute income are

54,000 and 16,000 respectively,

with Z's for normalisation,

a value of ₹ 73,600 for a income is transformed to?

$$\frac{73600 - 54000}{16000}$$

$$= \frac{19600}{16000} = 1.225$$

## Normalisation by Decimal Scaling

It is normalised by moving the decimal point of values of attribute A.

The no. of decimal point move, depends on the maximum absolute value of a.

A value  $v$  of  $A$ , is normalised to  $v'$  by comparing

$$v' = \frac{v}{10^j}, \text{ where } j \text{ is the}$$

smallest integer, such that maximum  $v'$  is less than 1.

## Data Integration

For a categorical data (discrete), a co-relation relationship b/w two attributes A and B can be discussed by 'Chi-Square'  $\chi^2$ .

Suppose that, a  $2 \times 2$  contingency table for the data is given below:

	male	female	Total
fiction	250	200	950
non-fiction	50	1000	1050
Total	300	1200	1500

for numeric  
data

Find the co-relation value of the given table.

$$Y_{AB} = \sum_{i=1}^n (a_i - \bar{a})(b_i - \bar{b})$$

$$\frac{N \bar{a} \bar{b}}{N}$$

$$= \frac{\sum (a_i b_i) - N \bar{a} \bar{b}}{N}$$

where,  $N$  = no. of tuples

$a_i, b_i$  = represent values of  $a$  and  $b$  in the tables

$\bar{a}, \bar{b}$  = represent mean value of  $A$  and  $B$ .

$\sigma_A \sigma_B = \text{Std deviation of A \& B}$

$$[-1 \leq \rho_{AB} \leq +1]$$

① if  $\rho_{AB} > 0$  then, A \& B are truly co-related.

② if values = 0, then A and B are independent, and there is no-correlation between them.

③ if  $\rho_{AB} < 0$  the A \& B are ~~truly~~ negatively co-related.

(Not important  
part of  
paper)

### Note

Suppose A has  $n$  distinct values,

$a_1, a_2, \dots, a_n$  and B has

$b_1, b_2, \dots, b_m$

The data tuples described by A \& B can be shown as a contingency table as:

Let  $(A_i, B_j)$  denote the event that

attribute A takes on the value  $A_i$  &  
attribute B takes on value  $B_j$

where,

$$A = A_i \text{ and } B = B_j$$

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

where,

$o_{ij}$  = observed frequency

and

$e_{ij}$  = expected frequency.

$$e_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{N}$$

for all  $i = 1, 2, \dots, r$  and  $j = 1, 2, \dots, c$

if  $H_0$  is true with null hypothesis

then  $\chi^2$  follows chi-square distribution

the value of  $\chi^2$  with n degrees of freedom

is called by  $\chi^2$  distribution

117  
117  
117

215  
ACD  
- cryptor -  
- M.L =  
Wn

5) Charlie

Date \_\_\_\_\_  
Page No. \_\_\_\_\_

## Numerical on Contingency

M      F      Total

million	250	C <sub>11</sub>	200	C <sub>12</sub>	450
non-million	50	C <sub>21</sub>	100	C <sub>22</sub>	1050
total	300		1200		1500

C<sub>1F</sub> = Count(male) Count(female)

$$\begin{aligned} & \text{male} \\ & \text{female} \\ & = \frac{300 \times 450}{1500} \\ & = 90 \end{aligned}$$

$$\begin{aligned} & C_{Fmale} = \frac{100 \times 450}{1500} \\ & C_{Ffemale} = \frac{100 \times 1050}{1500} \\ & = 300 \end{aligned}$$

Non-father

$$\begin{aligned} & C_{Mmale} = \frac{300 \times 1050}{1500} \\ & = 210 \end{aligned}$$

90    300  
210    840

$$C_{Ffemale} = \frac{1200 \times 1050}{1500}$$

$$\begin{aligned} & C_{Ffemale} = \\ & = 840 \end{aligned}$$

$$x^2 = \frac{(250-90)^2}{90} + \frac{(200-30)^2}{300}$$

$$+ \frac{(50-210)^2}{210} + \frac{(1000-890)^2}{840}$$

$$= \frac{(160)^2}{90} + \frac{(-160)^2}{360} + \frac{(160)^2}{210}$$

$$+ \frac{(160)^2}{840}$$

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

$$= 284.44 + 71.11 + 121.90 \\ + 30.97$$

$$= 507.92$$

$$= 22537$$

15-oct-18

53

Charlie

Date

Page No.

## Assignment - 2

Write short note on  
Data ~~structure~~ <sup>use</sup> of aggregation,  
Data compression  
Numerosity reduction  
Clustering  
Concept hierarchy generation.