

What is data warehouse:-

→ Data warehousing is a process of collecting and organizing large and complex data from various sources into a central repository, known as data warehouse, for reporting and analysis.

- The primary goal of data warehousing is to support decision making by providing business users with single, unified and integrated view of the data required to support planning and decision-making (e.g., financial analysis, operational decisions).
- A data warehouse is non-subject-oriented, integrated, time variant and non-volatile collection of data.

Subject-oriented:

A data warehouse is organized around major subjects such as customer, sales products, suppliers, etc. Rather than concentrating on day-to-day transactions of an organization, data warehouse focuses on modelling and analysis of data for decision makers.

Integrated:

A data warehouse is usually constructed by integrating various multiple heterogeneous sources such as relational database, flat files etc... and data cleaning and integration ensure consistency in data warehousing.

Time variant:

Data is stored to provide information from an historic perspective (past 5-10 years). Every key structure in data warehouse contains, either implicitly or explicitly a time element.

Non-Volatile

operational updates of data does not occur in data warehouse env. It does not req. transaction recovery processing and concurrency control mechanisms. It only requires initial loading of data and access of data.

## Difference between DBMS and data warehouse

- A database management system (DBMS) and a Data warehouse (DW) are both systems used to store and manage data, but they serve different purposes and have different characteristics.
- A DBMS is a software system that is designed to manage a database and provides users with the ability to interact with the database through various interfaces, such as SQL. The primary purpose of DBMS is to provide a centralized location for data storage and retrieval, and to enforce rules for data integrity, consistency and security.
- A data warehouse is a large, centralized repository of data designed for use by decision makers and business analysts. The data in a data warehouse is usually integrated from multiple sources, and it is transformed and summarized to support business intelligence activities such as data analysis and reporting. Unlike a DBMS a data warehouse is optimized for querying and analysis, rather than transactional processing.
- In summary, a DBMS is a general-purpose system for managing data, while a data warehouse is a specialized system for storing and analyzing large amount of historical data.

Difference b/w operational database system and Data warehousing.

OLTP (online Transaction Processing)

OLAP (online Analytical Processing)

- used by clerk, IT Professional (DBA)
- day-to-day operations
- Based on ER model
- DB design is application oriented
- DB size is 100MB - GB
- similar to DBMS
- read/write / index / prime key.
- purpose is to control & run fundamental tasks
- processing speed is very fast
- detailed view of data
- primitive and detailed data

database  
adminstrator

knowledge worker (eg: manager, analyst)

long-term informational requirements

decision support.

based on schema.

subject oriented.

DB size is 100GB - TB

similar to data warehousing.

lots of scans

Purpose is to help with planning, problem solving & decision support.

depends on the amount of data involved.

multidimensional view of data.

summarized data.

data

batched  
processes  
and  
queries

OLTP

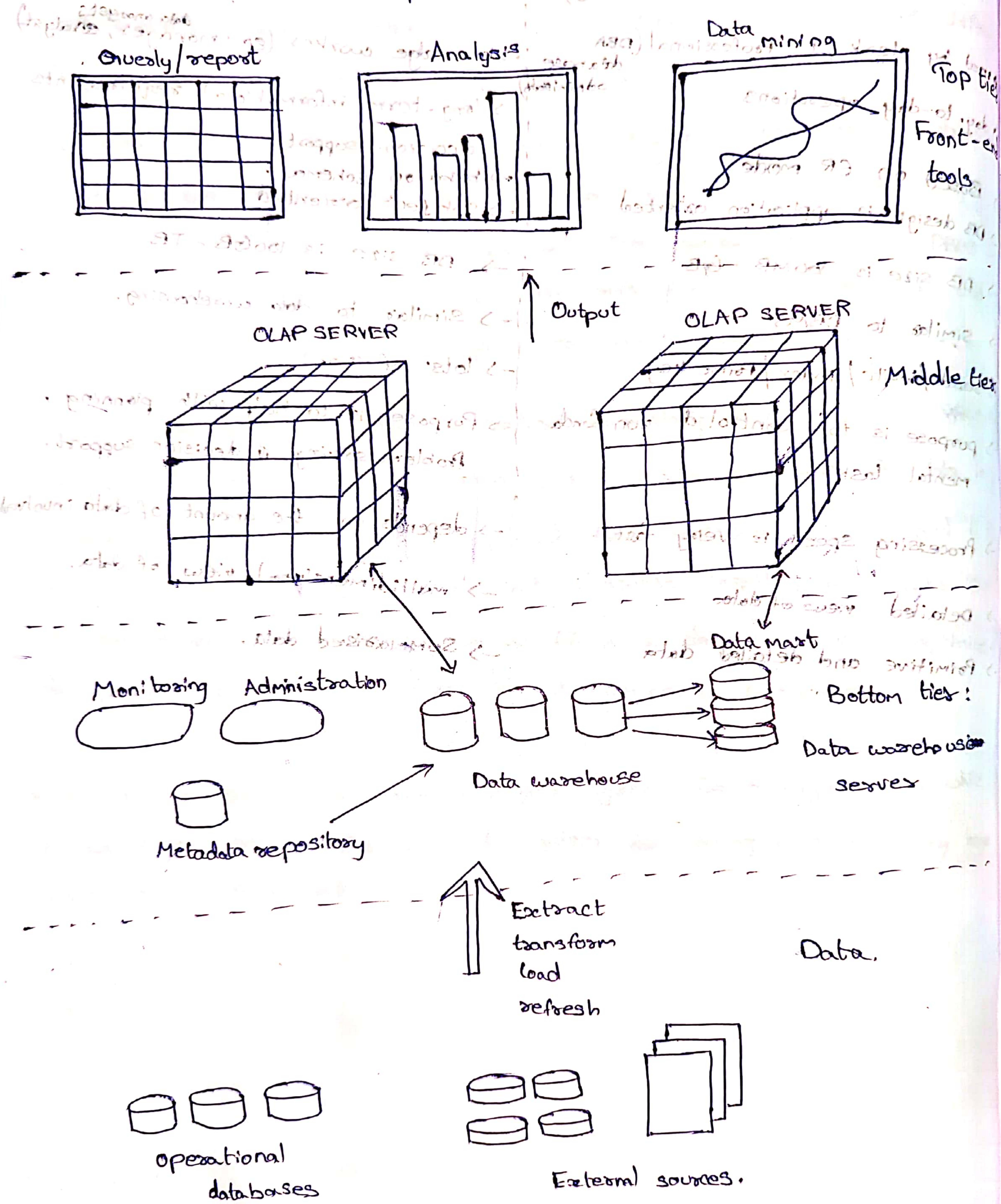
real time

OLAP

batched

### 3-Tier Architecture:-

→ Data warehouses often adopt a three-tier architecture.



Bottom tier: This is also known as Database Tier. This is the bottom-most layer and is responsible for storing the data. It communicates with the application tier to provide data when requested and to update the data when changes are made. This layer is also known as the data storage layer or the data layer.

Middle tier: This is also known as Application Tier. This layer is responsible for processing the data and handling all of the data flow (possible for multiple feeds) information. It is responsible for retrieving and business logic. It communicates with the database tier to retrieve and update the data and with the presentation tier to display the data. This layer is also known as the logical layer or the middleware layer.

Top tier: This is also known as presentation Tier. This is the top-most layer and is responsible for presenting data to the user. It usually contains the user interface and any client-side logic that is required to interact with the user. This layer is also known as the client layer or the user interface layer.

The three-tier architecture allows for separation of concerns, scalability, and improved security, as each tier can be managed individually and updated independently of the others. The 3-tier architecture is commonly used in web applications, where the presentation tier runs in a web browser, the application tier is hosted on a web server, and the database tier is hosted on a separate database server.

This diagram illustrates the three-tier architecture. At the top, there is a client application icon. Below it, there is a web server icon. At the bottom, there is a database server icon. Arrows point from the client application to the web server, and from the web server to the database server, indicating the flow of data between the three layers.

\* Data modeling: Process of designing schema of detailed and summarized info of datawarehouse. The main objective of the data modeling is to improve efficiency of the system and to support complex queries.

→ 3 levels :-

1. conceptual : mainly explains the semantics of data i.e simply the meaning of the data.
2. Logical : defines all the information (or) data in a structural format using data structures ; rules ; processes ; relationships ; etc ..
3. Physical : defines how the data is represented like in form the data is represented like tables (or) views etc..

\* Data warehouse models:- (Types)

1. Enterprise Warehouse :-

it will contains all the information about subjects related to a "entire" organisation both detailed and summarized . This is implemented on mainframes and super servers . Range from (gb to Tb)

2. Data mart :-

it will have data specific to group of users (not entire organisation) and only have summarized data.

3. Virtual warehouse :-

Containing the data copied from multiple source during the process of analysis and no separate schema or no base table.

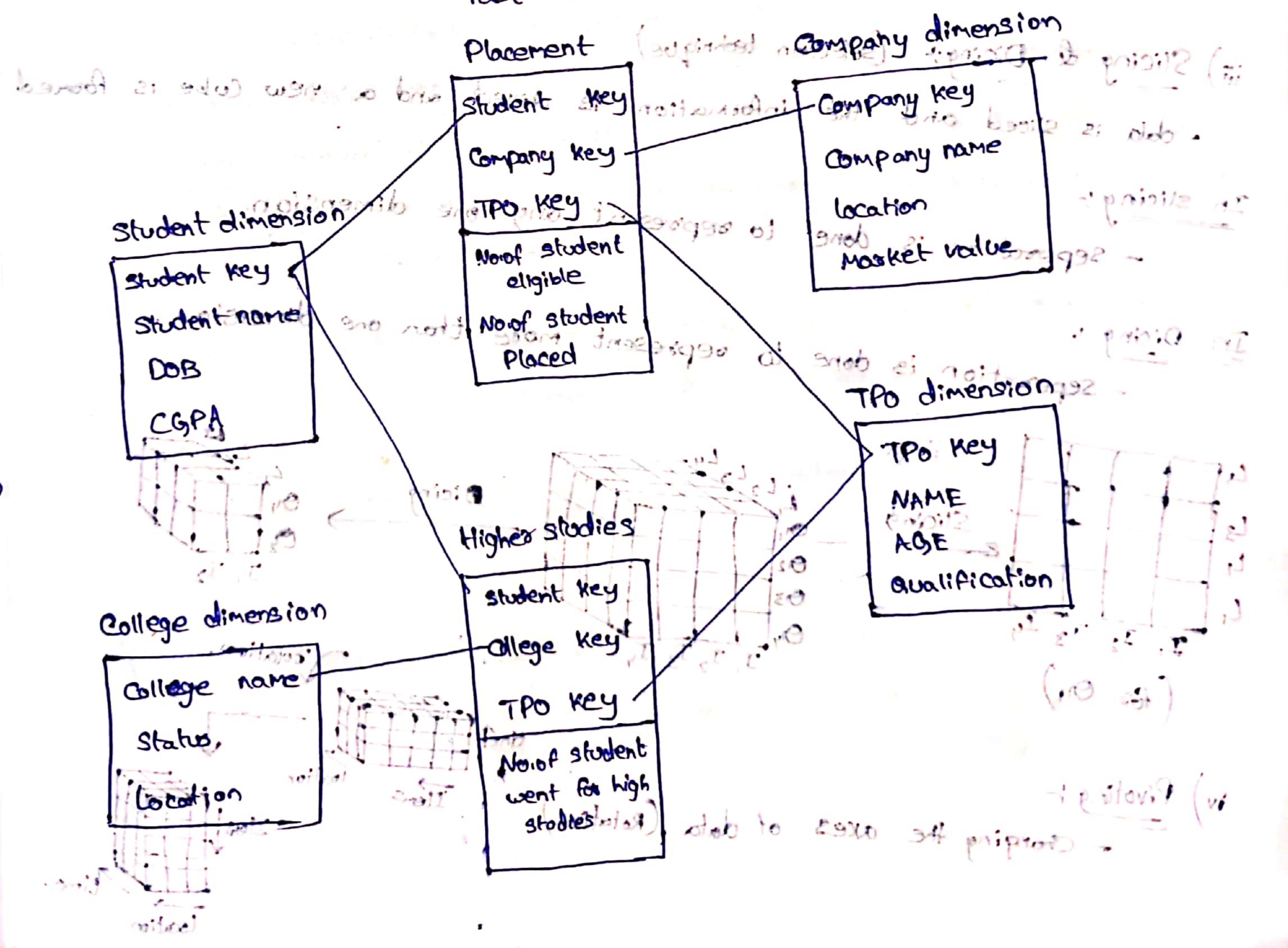
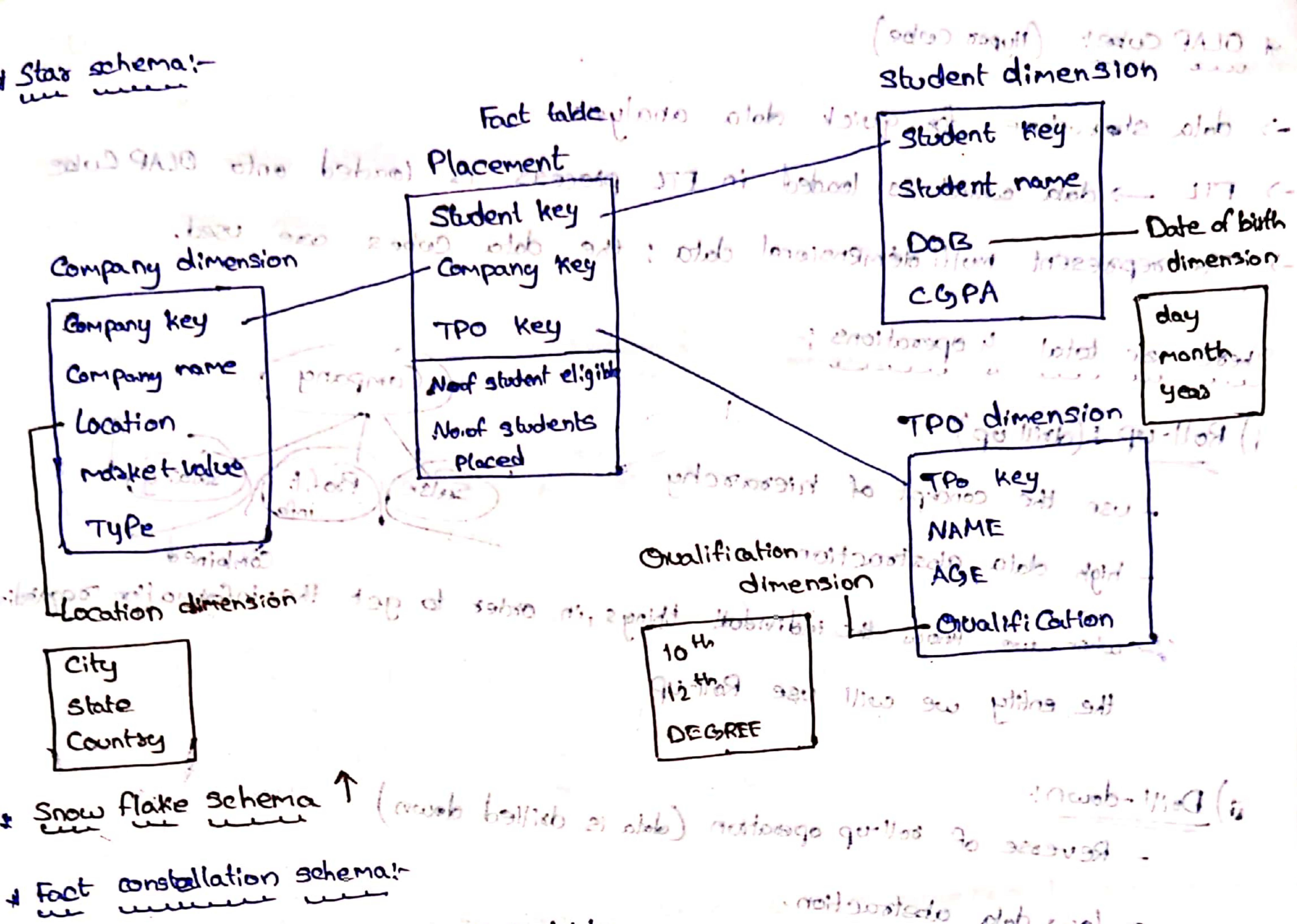
\* Fact Table:-

→ Relations between multidimensional data

\* Dimension Table :-

→ Tables related to each dimension and helps in describing dimension further.

\* Measures in a schema refer to numerical or quantitative data that is stored in schema



- \* OLAP Cube:- (Hyper cube)
- data structure for quick data analysis
- ETL → data which is loaded in ETL process is loaded onto OLAP Cube
- To represent multi dimensional data; the data cubes are used.

There are total 4 operations :-

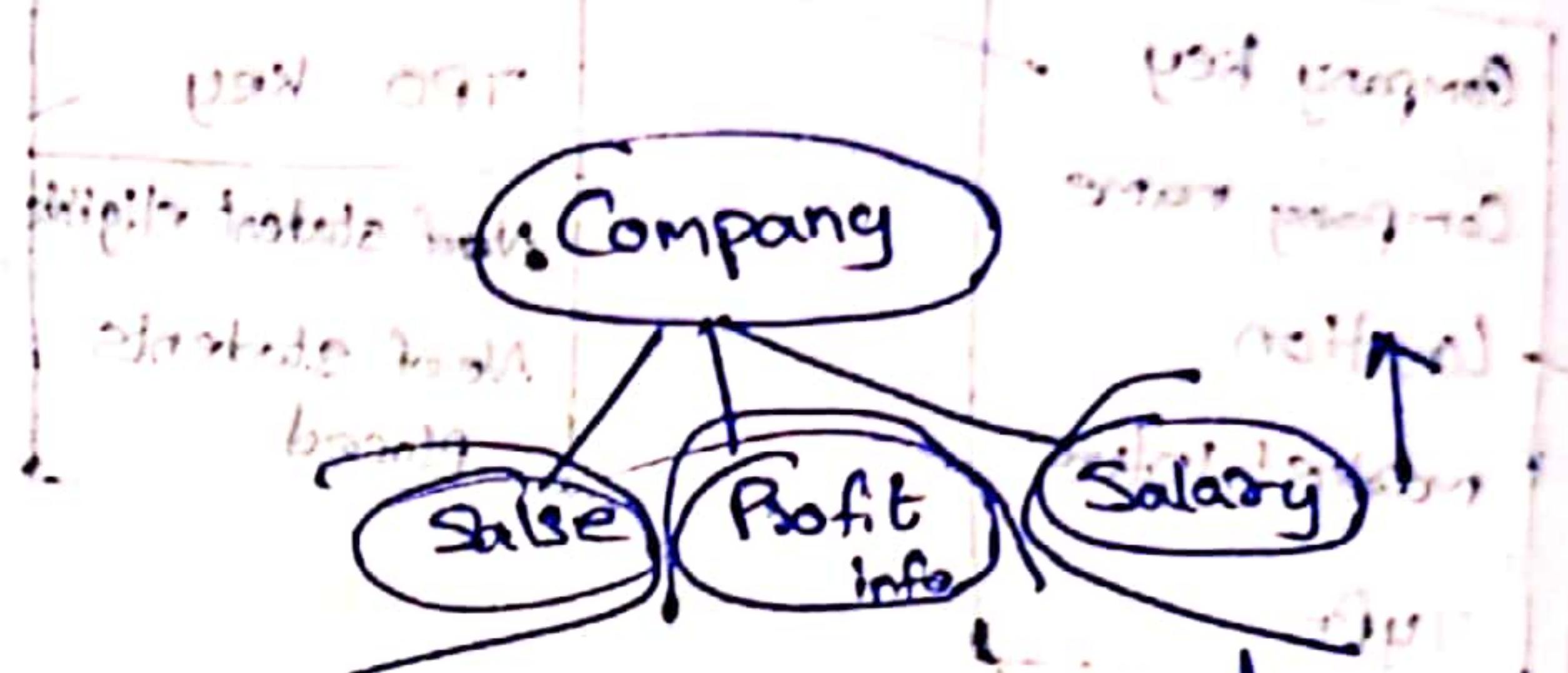
### i) Roll-up : (drill up)

- use the concept of hierarchy

- high data abstraction

- when we know the individual things, in order to get the information together

the entity we will use Roll-up



### ii) Drill-down:

- Reverse of roll-up operation (data is drilled down)

- low data abstraction.

### iii) Slicing & Dicing : (selection technique)

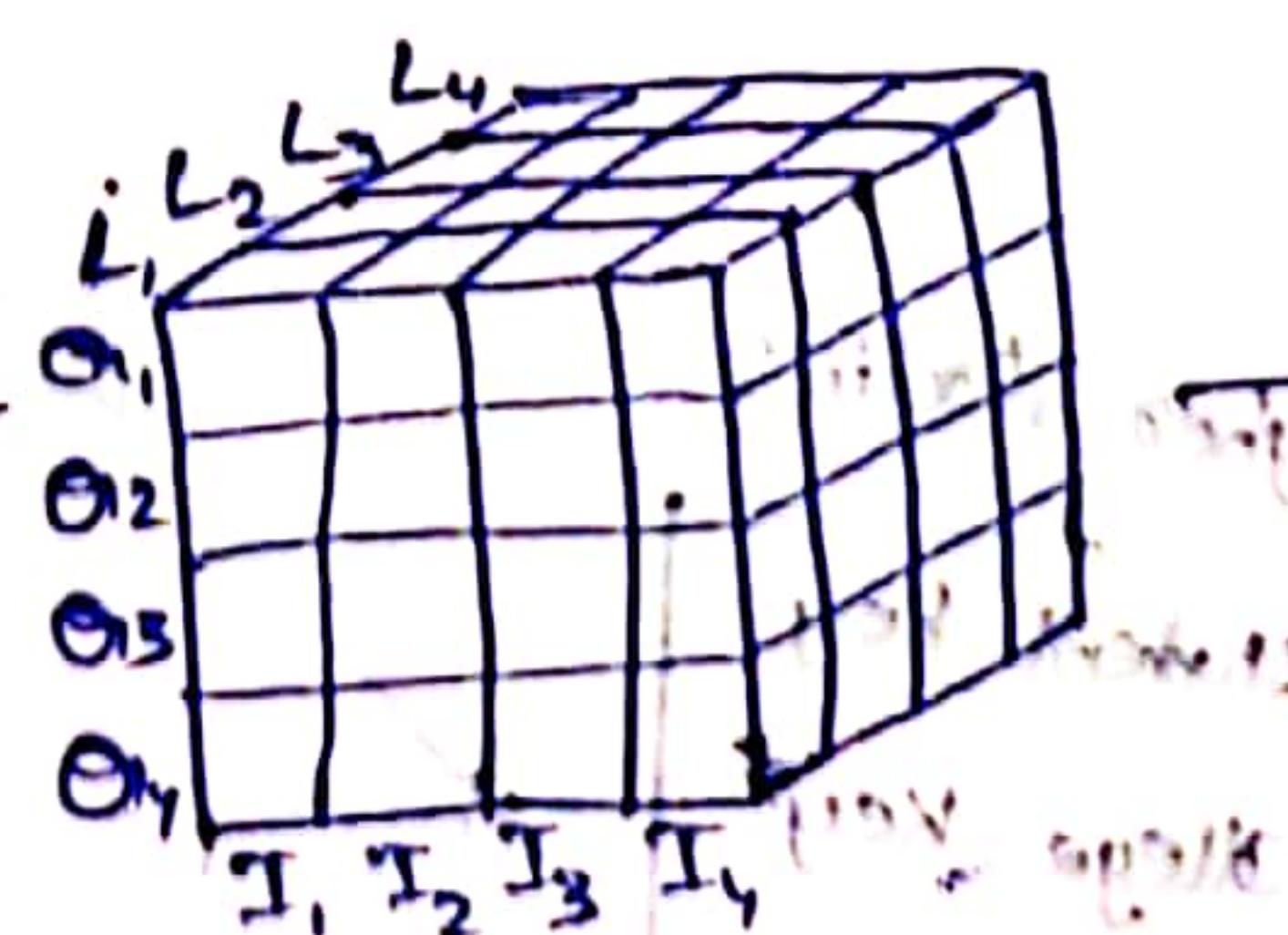
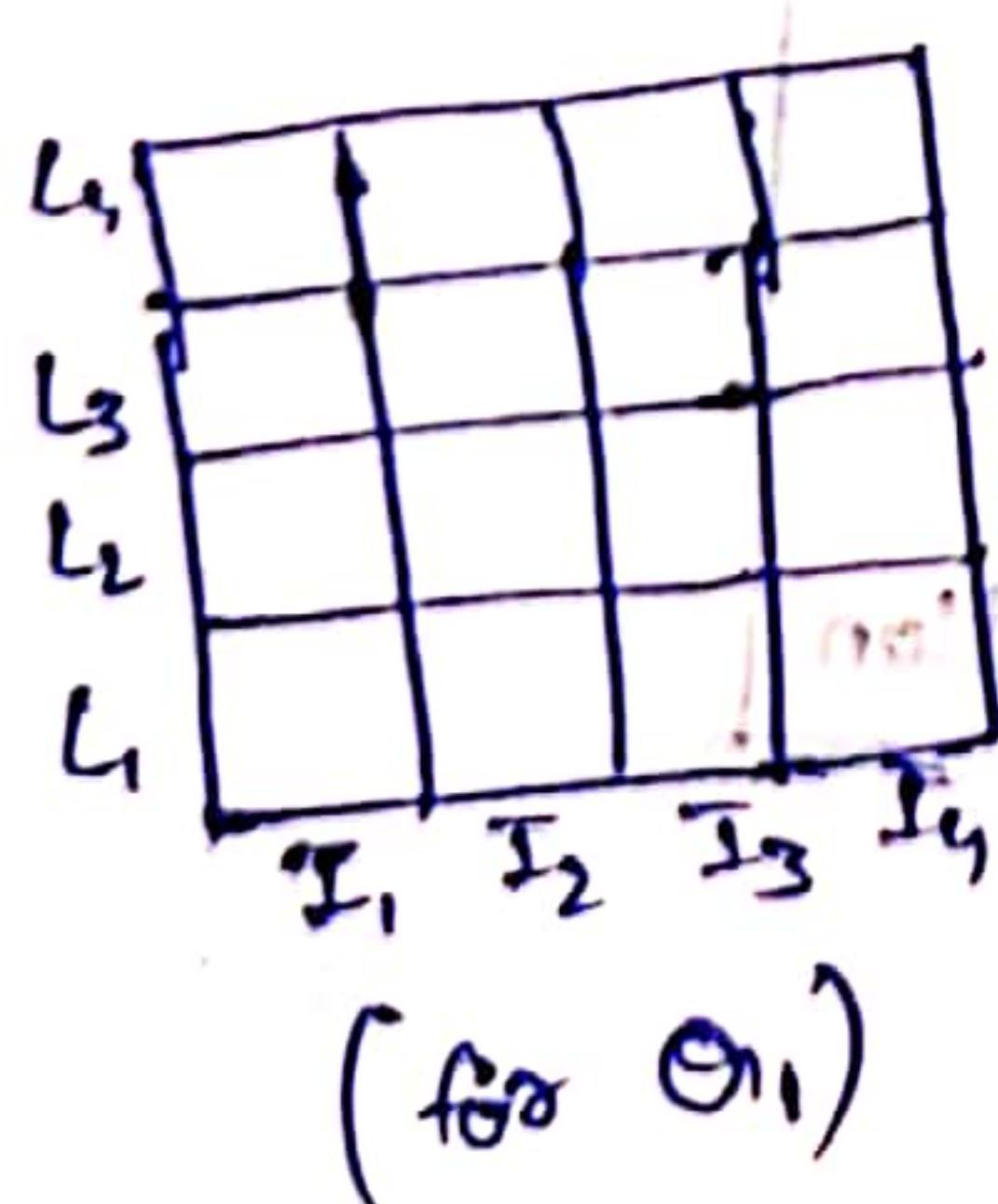
- data is sliced and the information is dived and a new cube is formed

In slicing :-

- separation is done to represent any one dimension.

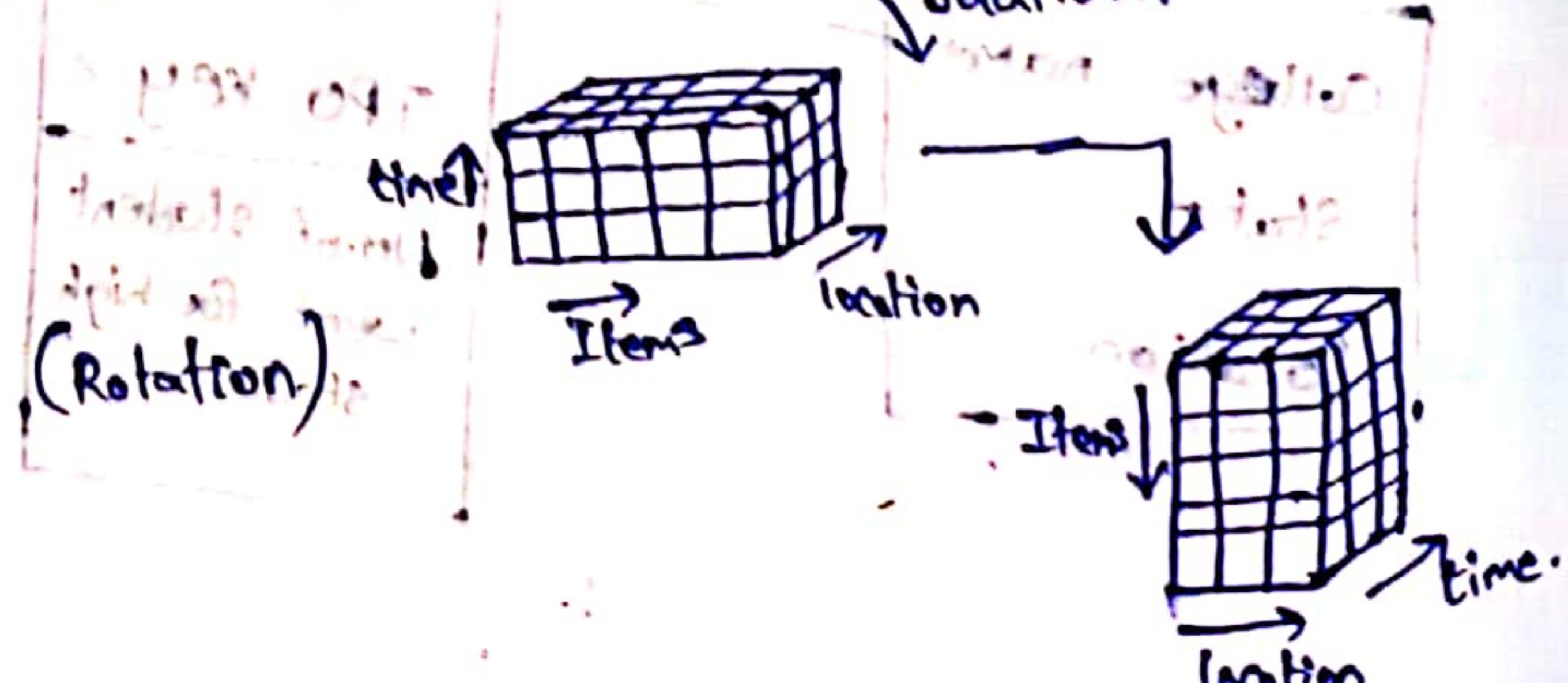
In Dicing :-

- separation is done to represent more than one dimension.

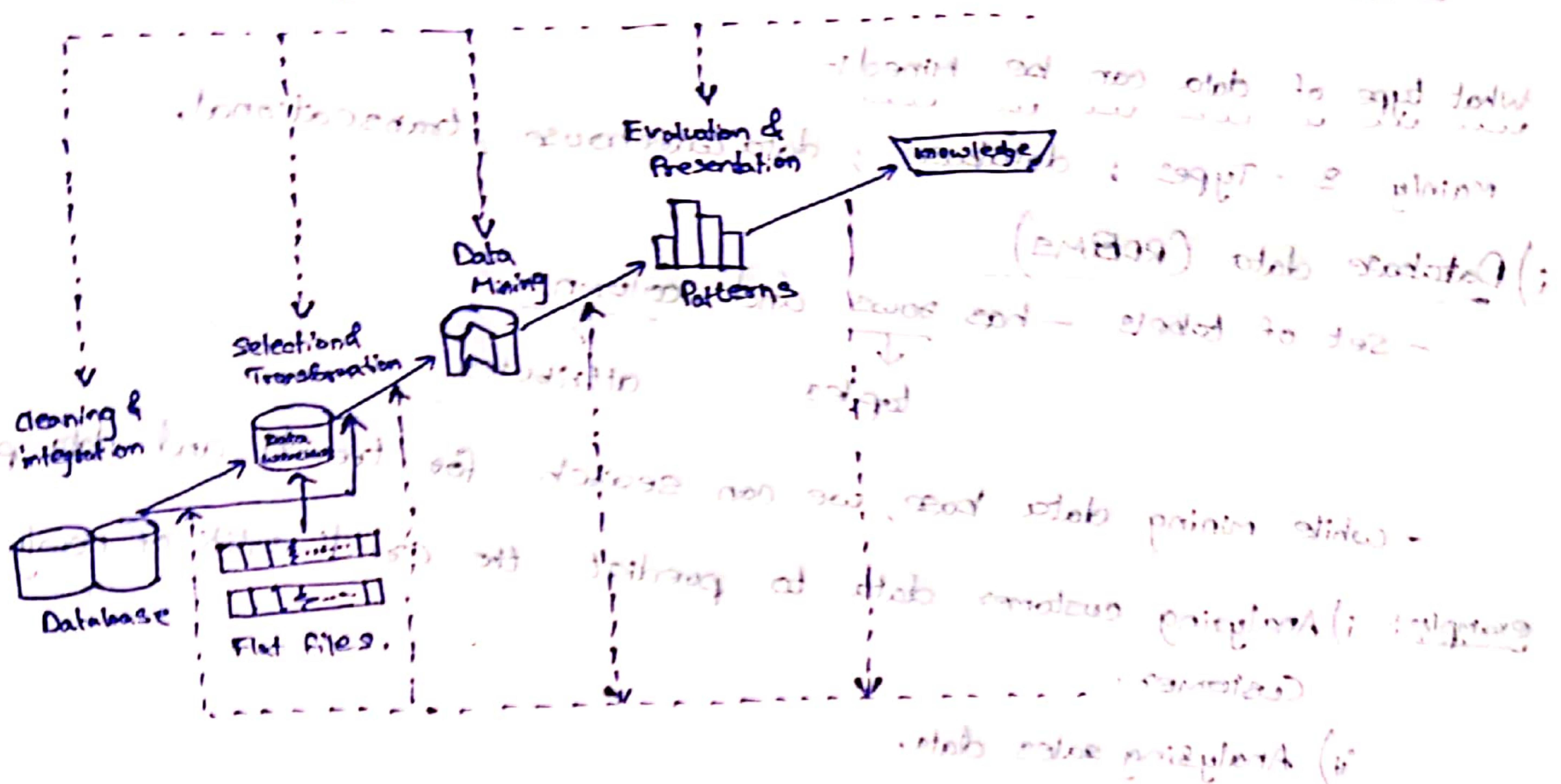


### iv) Pivoting :

- changing the axes of data, (Rotation)



**KDD**:  
The knowledge discovery in database process is a widely used methodology for extracting useful knowledge from the data. The KDD process typically involves the following steps:-



**Process Flow**

1. **Data Cleaning**: In this the raw data is cleaned and pre-processed to remove the inconsistencies, errors and missing values.

2. **Data Integration**: Data from multiple sources are combined into a single dataset.

3. **Data Selection**: Relevant data is filtered or chosen based on some criteria.

4. **Data Transformation**: The selected data is transformed into form by reducing that can be analyzed.

5. **Data Mining**: The transformed data is analyzed by using data mining algorithms to identify patterns, trends and relationships.

6. **Pattern Evaluation**: The identified patterns are evaluated to determine their significance and usefulness.

7. **Knowledge presentation**: The evaluated patterns are presented in a form that can be easily understood and interpreted.

## \* DATA MINING:-

- Data mining is defined as procedure of extracting info from huge sets of data.
- also defined as mining knowledge from data.

What type of data can be mined :-

mainly 3 - Types : database ; data warehouse ; transactional.

### i) Database data (RDBMS)

- set of tables - has rows and columns  
↓  
tuples      attributes
- while mining data base, we can search for trends and data patterns

example:- i) Analysing customer data to predict the credit risk of newly coming customer.  
ii) Analysing sales data.

### ii) Datawarehouse data:-

- Datawarehouse is a process of collection, and organizing the complex data from various sources into a central repository for reporting & analysis.
- In datawarehouse, data is stored in multidimensional structure (data cube), where each dimension will represent each attribute.

### iii) Transactional data:-

- Each record is called as transaction.
- Sales, flight booking, user clicks on web page.
- From this we can mine frequent patterns.

### Other-types of data:-

- Sequence data, data streams, spatial data, engineering design data, hypertext, multimedia, web data, etc.

\* Interestness of patterns:-

- This rise three questions :

i) what makes pattern interesting?

- easily understood by humans
- valid on new/test data
- potentially useful. (fulfill all needs)

ii) Can data mining system generate all of the interesting patterns?

- refers to completeness of a dm system.
- In reality it is impossible to generate all interesting patterns.

iii) Can data mining system generate only interesting patterns?

- refers to optimization of a dm system.
- generating only interesting patterns → challenging.
- if it is done, then it become easy and efficient for user.

\* Outlier analysis is a technique used in data analysis to identify and examine data points that are significantly different from majority of data points. these points lie far from center of distribution and may effect skewing results; decreasing accuracy or creating bias.

→ Outlier analysis involve:

- Identify the outliers (examin the data to find the outliers)
- analyzing the outliers (examin the outliers is it true or it is usual data point)
- Remove (or) handling outliers (remove or transform or use robust statistical methods)  
like, Z score; box plot; scatter plot, etc.

→ Outlier analysis ensure accuracy and reliability of results obtained from the statistical analysis.

→ Proximity measure or similarity (or) dissimilarity:

- Can take two or more states
- More than one state. There is implicit assumption that there is no relation between states.

$$d_{ij} = \frac{P-m}{P}$$

: where P refers to total number of attributes

& m refers to total patches.

Gender is symmetric attributes & remaining three are asymmetric attributes

Eg:-

1. Red A P<sub>1</sub>
2. Green B P<sub>2</sub>
3. Blue B P<sub>3</sub>
4. Red A P<sub>4</sub>

$$d(1,2) \Rightarrow \frac{3-2}{3} = \frac{1}{3} = 0.3$$

$$d(1,3) \Rightarrow \frac{3-0}{3} = 1$$

$$d(1,4) \Rightarrow \frac{3-2}{3} = \frac{1}{3} = 0.3$$

$$d(2,3) = \frac{3-1}{3} = \frac{2}{3} = 0.6$$

$$d(2,4) = \frac{3-0}{3} = 1$$

for Jack & Mary

$$d(3,4) = \frac{3-0}{3} = 1$$

$$t(Jack, Mary) = 3$$

$$P_i$$

$$P_j$$

$$N_i$$

$$N_j$$

$$\text{for gender} = \frac{x+s}{x+s+t}$$

$$r(Jack, Mary) = 2$$

$$s(Jack, Mary) = 0$$

$$t(Jack, Mary) = 1$$

$$d(Jack, Mary) = \frac{0+1}{2+0+1} = \frac{1}{3} = 0.3$$

Jack & Jim :-  $\alpha = 1$ ;  $\gamma = 1$ ;  $s = 1$ ;  $t = 3$ .

for Jack & Jim :-  $\frac{1+1}{1+1+1+0} = \frac{2}{4} = 0.5$

Many & Jim :-  $\alpha = 1$ ;  $\gamma = 2$ ;  $s = 1$ ;  $t = 2$   
 $d(M, J) = \frac{\alpha+s}{\alpha+\gamma+s+t}$

for Many & Jim :-  $\frac{1+2}{1+2+1+0} = \frac{3}{4} = 0.75$

- Proximity measure for binary attributes is 1 minus of attributes i.e. equal after both if attributes are equal to 1. And for similar :-
- Proximity measure for cardinal attributes.

distance measure for asymmetric binary attributes

$$d(i,j) = \frac{\alpha + s}{\alpha + \gamma + s + t}$$

→ Proximity measure similarity (or) dissimilarity.

In simplest words, it is the degree of similarity & dissimilarity between two or more states.

→ Can take two or more states.

States need not match. None or complete matches also.

$$d_{ij} = \frac{P-m}{P} ; \text{ where } P \text{ refers to total number of attributes}$$

& m refers to total matches.

Eg:-

1. Red A P<sub>1</sub>

$$d(1,2) \Rightarrow \frac{3-2}{3} = \frac{1}{3} = 0.3$$

2. Green B P<sub>2</sub>

$$d(1,3) \Rightarrow \frac{3-0}{3} = 1$$

3. Blue B P<sub>3</sub>

$$d(1,4) \Rightarrow \frac{3-2}{3} = \frac{1}{3} = 0.3$$

4. Red A P<sub>4</sub>

$$d(2,3) = \frac{3-1}{3} = \frac{2}{3} = 0.6$$

$$d(2,4) = \frac{3-0}{3} = 1$$

$$d(3,4) = \frac{3-0}{3} = 1$$

→ Proximity measure for ordinal attributes.

→ Proximity measure for binary attributes.

$$d_{ij} = \frac{s+t}{q+r+s+t}$$

where, q = no. of attributes i.e. equal to both i & j.

s = no. of attributes i.e. equal to "0" for i & "1" for j.

t = no. of attributes i.e. equal to "0" for both i & j.

distance measure for asymmetric binary attributes.

$$d_{ij} = \frac{s+t}{q+r+s+t}$$

Ex:-

Name	Gender	Fever	Cough	Test1	Test2	Test3	Test4
		No	P <sub>1</sub>	No	P <sub>1</sub>	No	P <sub>1</sub>
Jack	M	y	N <sub>o</sub>	P <sub>1</sub>	N <sub>o</sub>	N <sub>o</sub>	(i,i) b
Mary	F	y	N <sub>o</sub>	P <sub>1</sub>	N <sub>o</sub>	N <sub>o</sub>	(N <sub>o</sub> , N <sub>o</sub> ) N <sub>o</sub> b
Jim	M	y	P <sub>1</sub>	N <sub>o</sub>	N <sub>o</sub>	N <sub>o</sub>	b

Gender is symmetric attribute & remaining three are asymmetric attributes

~~for gender =  $\frac{x+s}{x+s+t}$~~

~~$\Rightarrow q = q_{(Jack, Mary)} = 2 ; \gamma = 0 ; s_{(Jack, Mary)} = 1 ; t = 3$~~

~~$t_{(Jack, Mary)} = 3$~~

~~for Jack & Mary =  $\frac{0+1}{2+0+1} = \frac{1}{3} = 0.\bar{3}$~~

~~Jack & Jim :-  $q = 1 ; \gamma = 1 ; s = 1 ; t = 3$ .~~

~~for Jack & Jim :-  $\frac{1+1}{1+1+1} = \frac{2}{3} = 0.\bar{6}$~~

~~Mary & Jim :-  $q = 1 ; \gamma = 2 ; s = 1 ; t = 2$~~

~~for Mary & Jim :-  $\frac{2+1}{1+2+1} = \frac{3}{4} = 0.75$~~

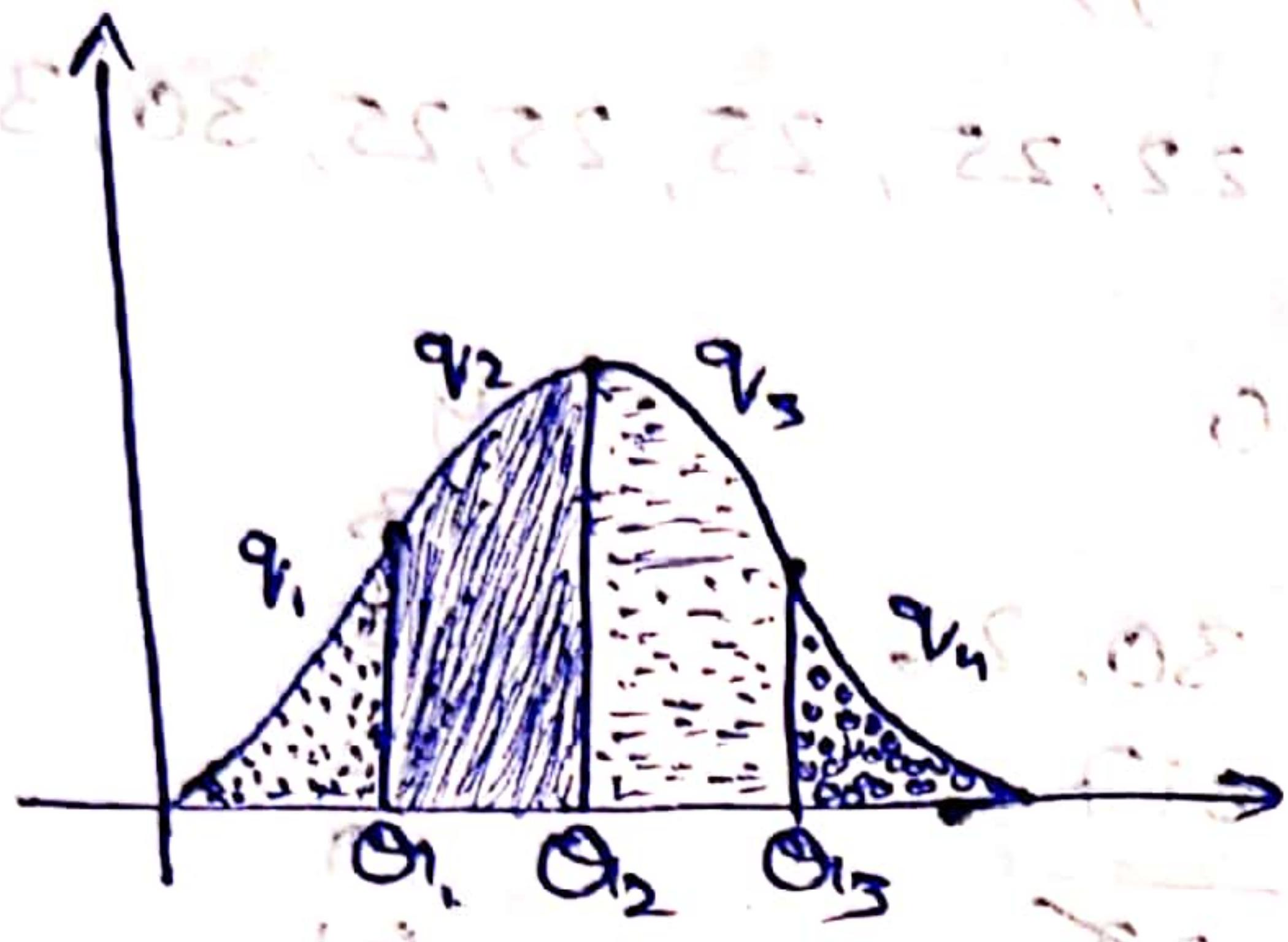
01/02/2023 | 18

\* Quantiles :- are lines which divide the graph into four equal parts.

i.e. (proportion 25%) original data set will be divided.

here,  $Q_1, Q_2, Q_3$  are Quantiles.

$Q_1, Q_2, Q_3, Q_4$  are quarters. (each of 25%)



Q) 13, 15, 16, 16, 19, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40,

45, 46, 52, 70.

A) As we know that in

$$\text{Mean} = 30$$

$$\text{Median} = 27.5$$

$$\text{Mode} = 25, 35$$

$$N = 26$$

$$Q_1 : 13, 15, 16, 16, 19, 20$$

$$Q_1 = 20$$

$$Q_2 : 21, 22, 22, 25, 25, 25$$

$$; Q_2 = 25.$$

$$Q_3 : 25, 30, 33, 33, 35, 35, 35$$

$$Q_3 = 35$$

$$Q_4 : 35, 36, 40, 45, 46, 52, 70$$

\* Box Plot:-

(Box Plot) shows distribution, center & spread of the data.

→ Ends of the box plot represent the quartiles ( $Q_1$  &  $Q_3$ )

→ The division in the center of the box plot is "median".

→ The two extreme ends outside of the box plot are smallest value &

largest value.

- These lines are known as whiskers.
- $\rightarrow$  The five point summary of a distribution consists of  $\Theta_1; \Theta_2; \Theta_3;$  (here  $\Theta_2 = \text{median}$ ) smallest value; largest value.

(Q) Suppose, that a hospital tested the age and body fat data for 18 randomly selected adults with the following results, first calculate the mean, median of age & fat. Draw the box plot for age & fat.

Age	23	23	27	27	39	41	47	49	50	52	54	54	56	57	58	58	60	61
fat%	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7

$N = 18$

Ans) Age :-  
mean :- 23  
23  
27  
7.8  
9.5  
17.8  
25.9  
26.5

$$\text{mean} = \frac{836}{18} = 46.4$$

$$\text{median} = \frac{50 + 52}{2} = 51$$

mode = 23, 27, 54, 58.

$$\Theta_1 = 39 ; \Theta_2 = 51 ; \Theta_3 = 57.$$

$$\text{mean} = \frac{518.1}{18} = 28.783$$

$$\text{median} = \frac{31.2 + 34.6}{2} = \frac{65.8}{2} = 32.9$$

mode =

$$\Theta_1 = 31.4 ; \Theta_2 = 32.9 ; \Theta_3 = 30.2$$

