

# Module-1

6/1/23  
Monday

## Data Mining

Extract of interesting (non-trivial, implicit, previously unknown and potentially useful) information or patterns from data ~~large~~ database. processed form of data is known as information.

KDD - knowledge discovery in database.

knowledge extraction.

data/pattern analysis.

data archaeology.

data dredging.

information harvesting.

business intelligence.

Non data mining.

(Deductive) query processing.

KDD is also known as datamining.

## Steps in KDD

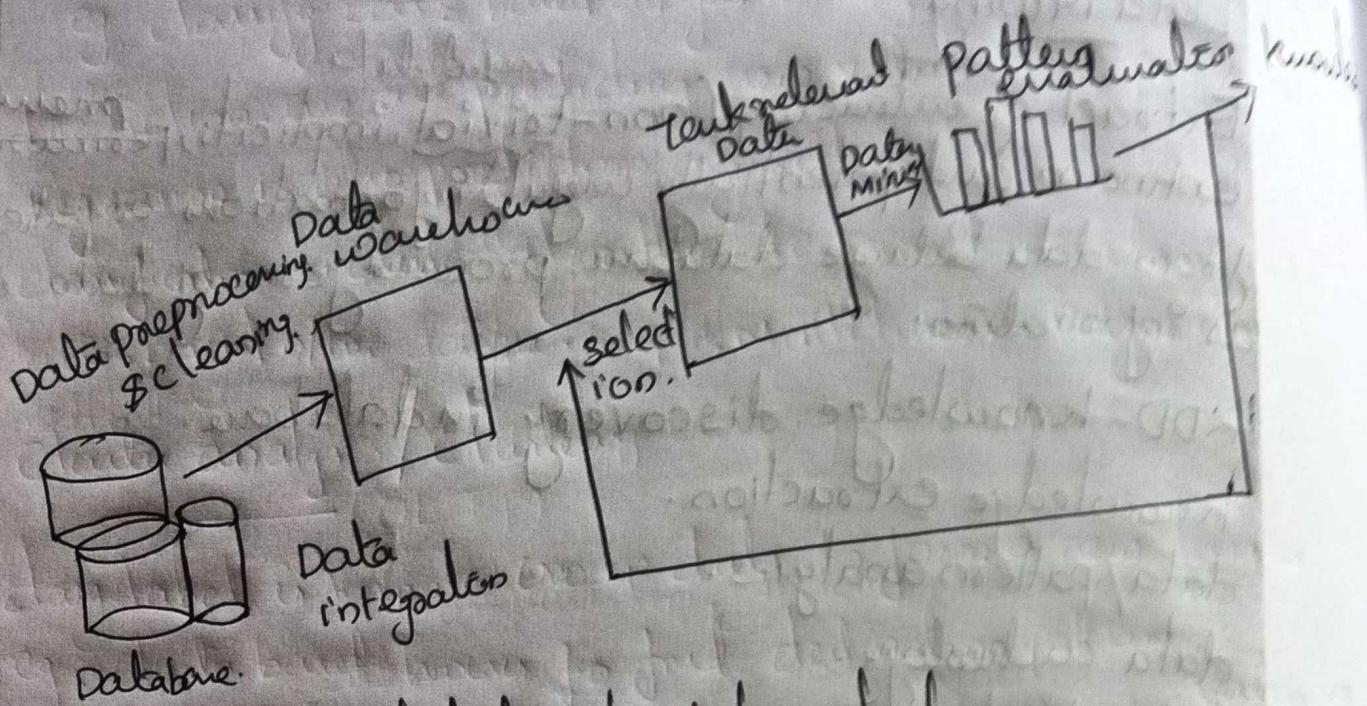
1. Data preprocessing and cleaning.

2. Data selection and Transformation.

3. Data Mining

4. Pattern evaluation.

Data mining: the core of knowledge discovery process.



Creating a target data set : data selection.

Data cleaning and preprocessing : removing unwanted data.

2) Find useful features dimensionality.

choosing functions of data mining

summarization, classification, regression, association, clustering

choosing the mining algorithm.

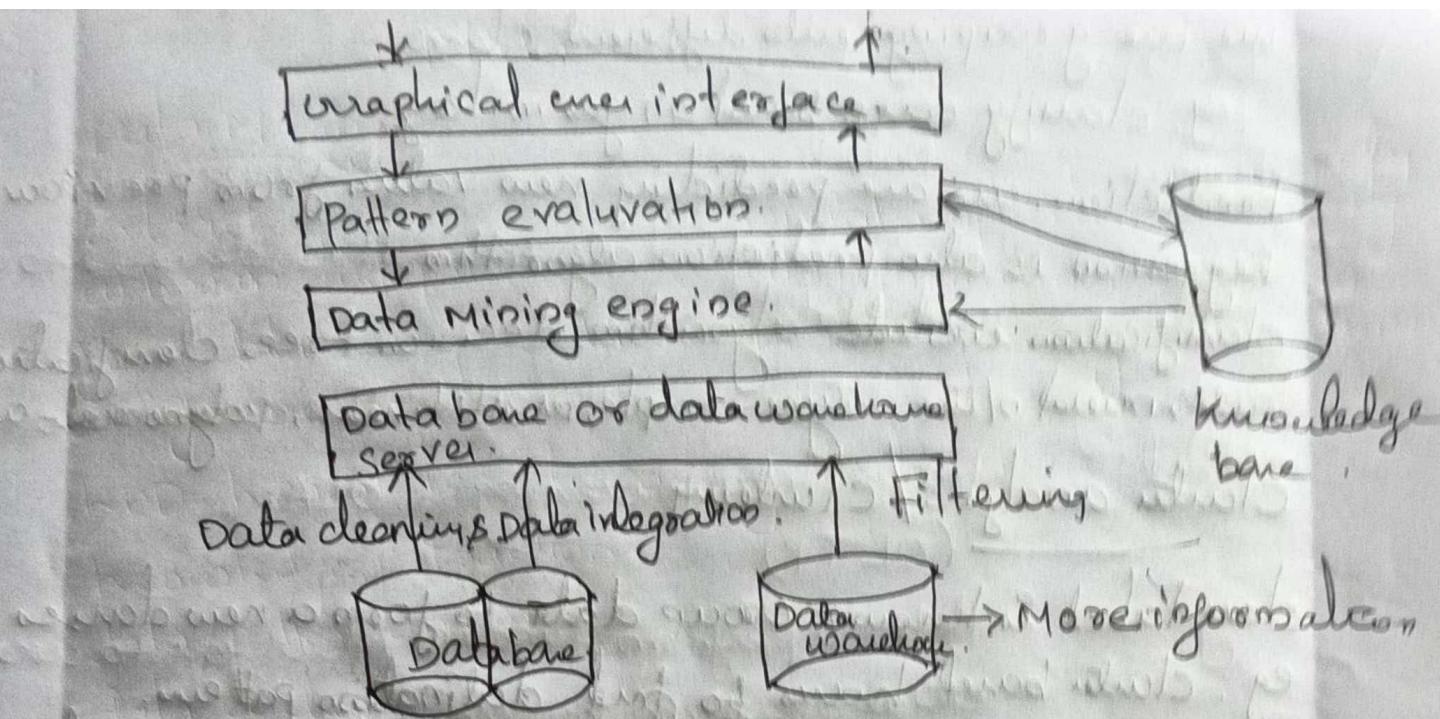
Data mining : search for patterns.

Pattern evaluation & knowledge

visualization, transformation, removing redundancy pattern.

use of discovered knowledge.

Architecture a Typical Data mining System.



### 6 Major Components

1. Database, data warehouse, www, or other information repository.
2. Database or data warehouse server.
3. knowledge base.
4. Data Mining engine
5. pattern Evaluation module
6. user interface.

### 12/22 Tuesday: Datamining functionalities

- Concept description: characterization and discrimination
- Generalize summarize and constant data characteristic eg: dry vs wet regions
- Association (correlation and causality)
  - Multidimensional vs single dimensional associations
- classification and prediction.
  - Finding models that distinguish class or concept for future predictions

we are grouping into different classes.

eg: classify cars based on gas mileage.

Prediction: we are predicting new values from previous values.  
Regression is also prediction algorithm.

Classification include algorithms, decision tree, classification rule, neural network type - simple, multi linear, polynomial - one layer

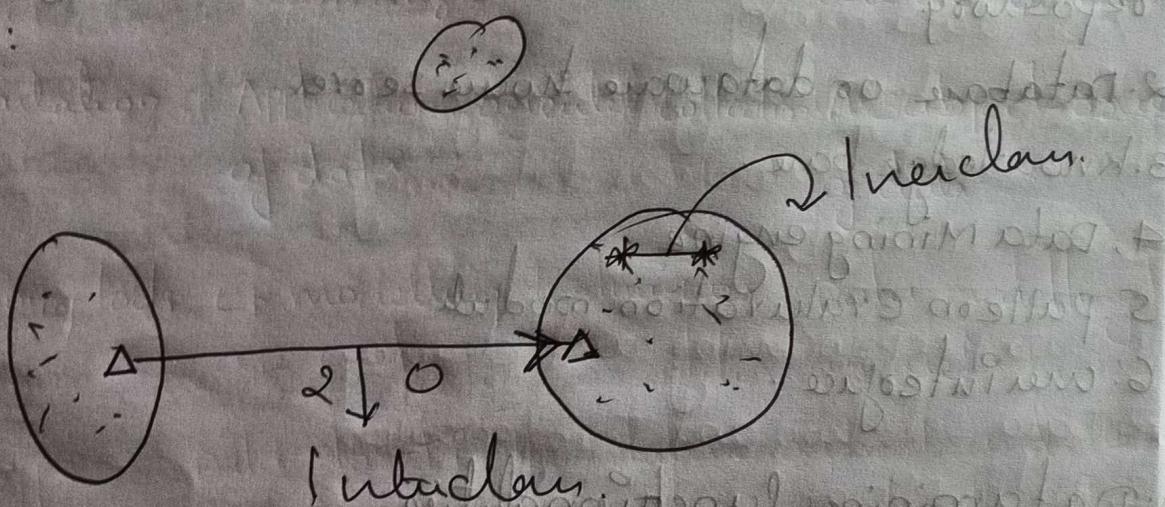
### cluster analysis / clustering

class label is known: group data to form a new classes.

eg: cluster based houses to find delimitation patterns.

clusters based on the principle maximize ~~b/w the class~~ the <sup>value ↑</sup> in b/w clusters and minimizing the <sup>value ↑</sup> interclass similarity.

eg:



similarity distance measure

Inter class similarity - b/w the class.

Intra " " - within the class.

### Outlier analysis (noise)

Outlier: A data object that does not comply with general behavior of the data.

Outlier Analysis: It can be considered as noise or exceptions but it is useful in fraud detection and error analysis.

## Trend and evaluation analysis

Trend and deviation: regression analysis similarly based analysis

Segmented pattern mining, periodicity analysis  
other patterns: directed or statistical analysis

## 7. Major issues in data mining

- Mining methodology and user interaction
- Mining different kinds of knowledge in database
- Performance and scalability
- Issues relating to the diversity of the datatypes.
- Issues related to application and social impacts.

Human interaction.  
query fitting  
changing data.

## 9/2/23 Thursday Data warehousing

A datawarehouse is subject-oriented, integrated, time-variant non-volatile collection of data in support of management's decision making process.

Decision support database.

Datawarehousing: The process of conducting and using data warehouse.

Time variant: entirely different from database.

non-volatile: operational update of data does not occur in the datawarehouse environment.

Initial loading of data source of data

## X. Database v/s Datawarehouse

Feature

Database

Datawarehouse

Purpose

It is designed to record data

It is designed to analyse data

Processing method.

The database uses the online transaction processing (OLTP) priority.

Usage

Allows you to analyze your business.

Tables & joins complex as they are normalized

Simple because they are denormalized.

Orientation

A subject-oriented collection of data.

Designing tools

ER modelling techniques are used.

Data type

Data stored in the database is up-to-date.

Current and historical data is stored in data warehouse.

May not be up-to-date.

Query type

Simple hierarchical queries are used.

Data summary. Detailed data stored in the database. It stores highly summarized data.

OLTP - Online Transaction processing.

OLTP

uses  
functions  
DB design  
Data

usage  
access

unit of work

# records of around

# users

DB size

metric

clerk, IT professional knowledge workers.  
day-to-day operation. decision support  
application oriented. subject oriented.  
current up to date detailed, historical, summarized  
flat relational isolated multidimensional, integrated  
repetitive consolidated  
read/write/index/hash ad-hoc  
on primary key. lots of scans.  
short, simple transaction complex query.  
items millions.

3/2/23 Tuesday A multidimensional data model  
Monday A multidimensional data model

A data warehouse is based on a multidimensional which view data in the form of a data cube. Data is represented in the form of data.

location = "Chicago" loc = "New York" loc = "Toronto" loc = "Vancouver"

item

item

item

item

home

home

home

home

time ent. comp. phone sec. ent. comp. phone sec. ent. comp. phone sec. ent. comp. phone sec.

Q1

Q2

Q3

Q4

OLAP

Cello

## Data Cube

— Data cube allows data to be modelled and viewed in multiple dimensions.  
It can be defined by dimensions and facts.

## Dimension

— Dimensions are the perspectives or entities, which is table, store entity, fact-numerical tables, relationship b/w dimensions tables.

## Example of dimension table:

time
time-key
day
day-of-the-week
month
quarter
year

branch
branch-key
branch-name
branch-type

eg of fact table

Problem

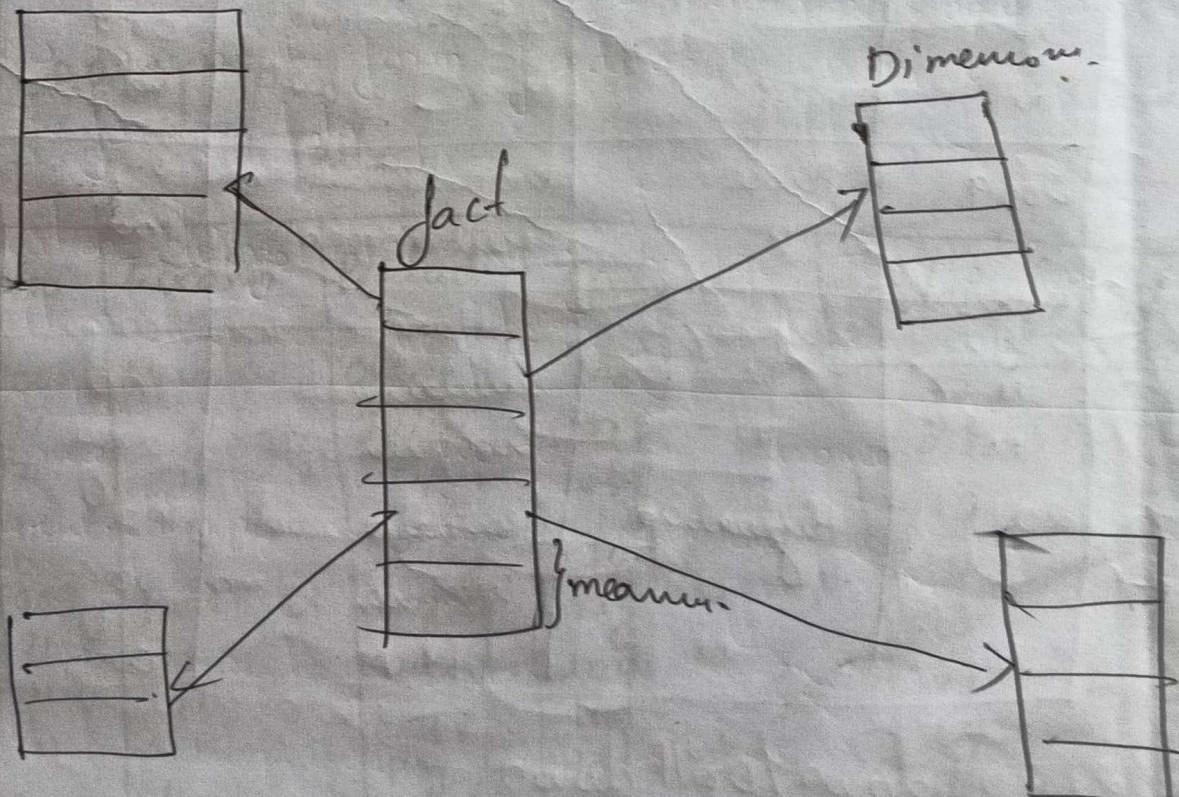
time-key.
item-key
branch-key.
location-key.
units-sold.
dollars-sold
avg-sales.

## X. Conceptual Modelling of Data warehouse / Data warehouse schema

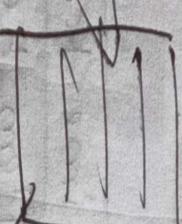
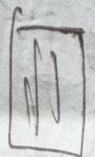
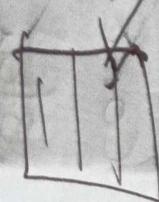
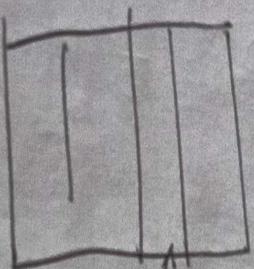
Star schema - fact table <sup>connected</sup> is the middle to a set of dimension table.

Snowflake schema : A refinement of star schema to normalize in a set of smaller dimension tables, forming a shape similar to.

Galaxy schema : Multiple fact tables share dimension tables, viewed as a collection of states through stars therefore called galaxy schema .



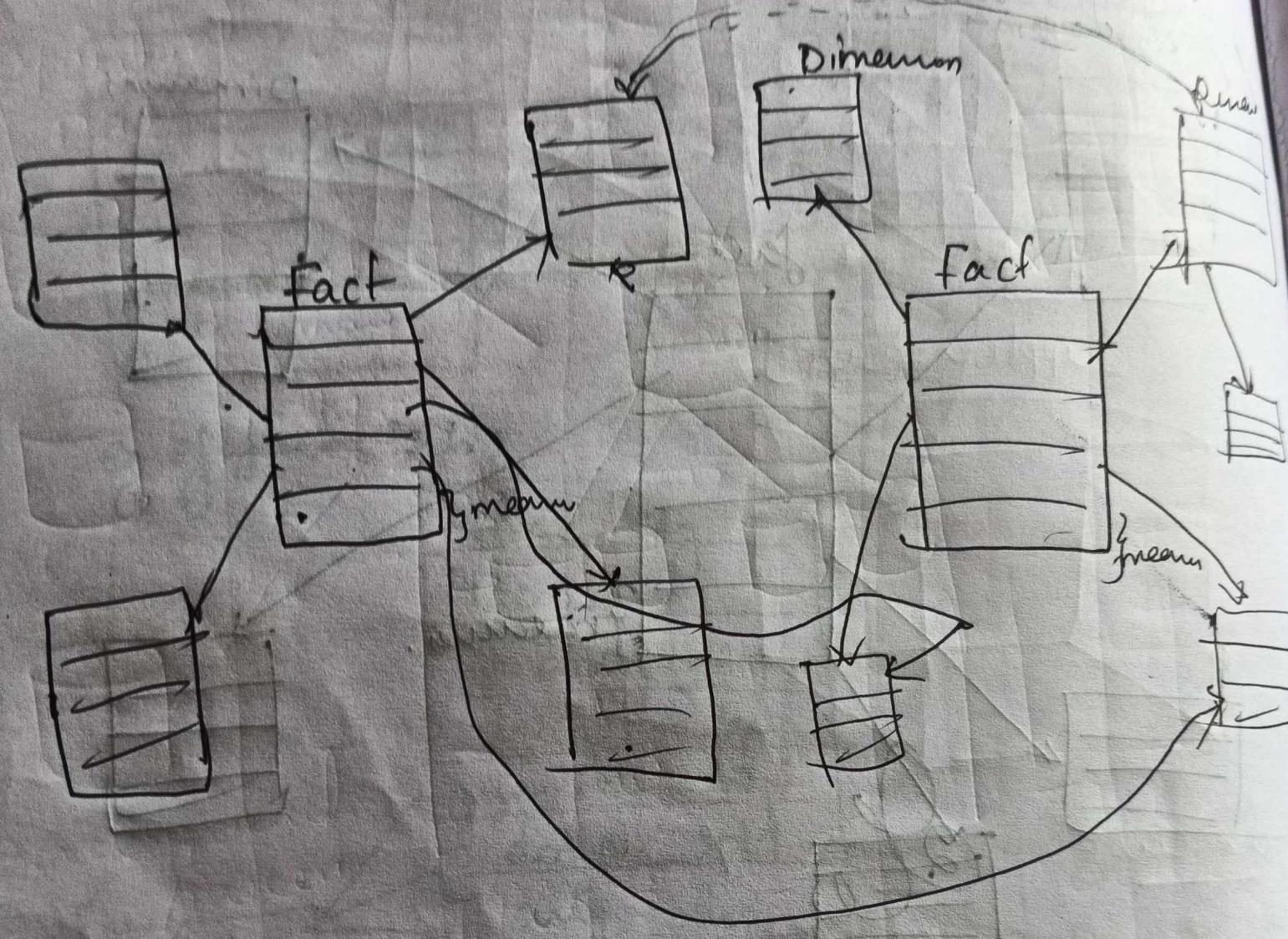
memory  
snowflake scheme



coloring scheme

coloring scheme  
coloring scheme

crabady schens



14 Feb 22  
Tuesday

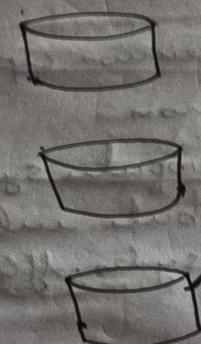
Datawarehouse architecture

14/2/22  
Tuesday

## Data ware house architecture

data warehouse architecture

operational data.

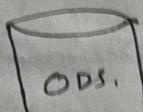


Background proc.

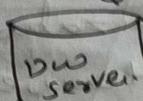
Load  
cleaning  
integrating  
transform

METADATA

Reconciled Data.



DRIVEN data



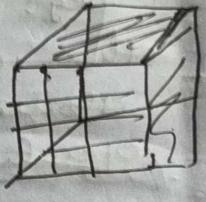
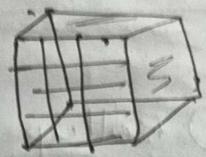
Analytic

Data marts



TIER-1  
DW server

Data mining  
application



TIER-2  
OLAP enabling  
OLAP server

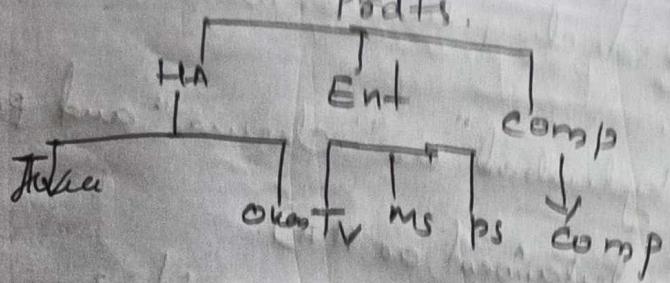
15/2/23  
Wednesday

## OLAP Operations

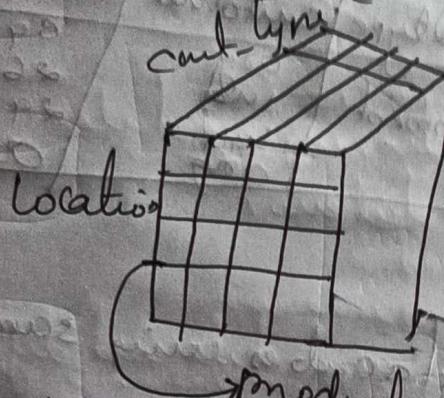
Table 1

Year	Location	Product	Customer	Total sales
2008	Mumbai	Overs	Ind	32
"	"	TV	Org	34
"	"	TV	Org	12
"	"	Municgym	Ind	11
2009	Pune	Computer	Org	43
"	"	Municgym	Ind	21
"	"	TV	Ind	14
"	"	Play_st	Ind	10
2010	Chennai	Overs	Org	9
"	"	Municgym	Ind	8
"	"	Computer	Org	13
				29

• used to decrement the dimension hierarchy by one level.

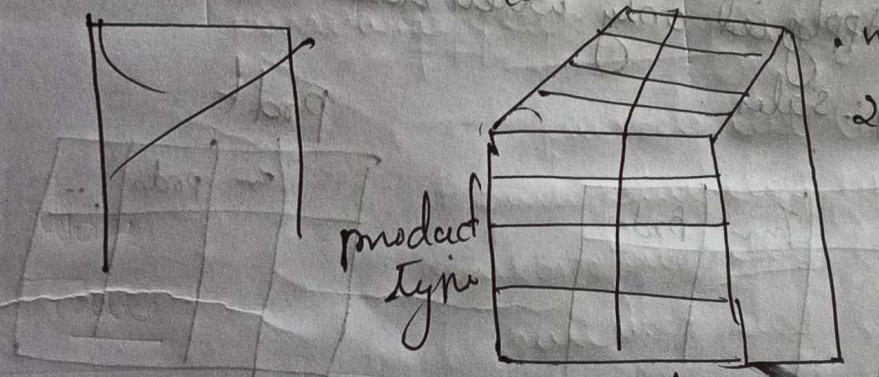


3) Pivot operation (rotate)



It is used to view data cubes in different form.

• view upto 2D planes.



4) slice : It is used to extract a slice of the original cube corresponding to a single value of a given dimension. (similar to select operation RDBMS)  
for eg.. slice Table 2 on location = "Hyderabad".

year	location	product-type	category	Total sales
2011	Hyderabad	entertainment	ind org	29
"	"	computer system	org.	18
"	"	TA	org.	03

16/2/23  
Thursday

### 5) Dice operation

It is used for selecting two or more dimensions and two or more values.

e.g.: location ("Hyderabad" or "Chennai" and product type ("Entertainment" or HA)).

After applying dice operation from Table 2.

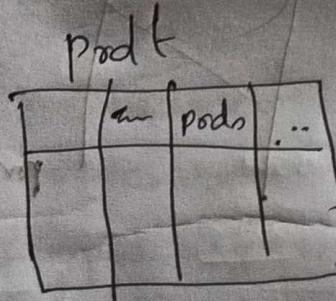
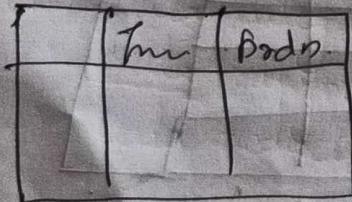
Year	location	product type	customer	Total sales.
2010	Chennai	HA	Org	89
"	Hyderabad	Entertainment	Ind	26
2011	"	HA	Org	29
"	"	"	Org	03

### 6) Drill Across

Analyse the cells of a data cube using same set of dimensions and move data from one data cube to another.

Applied only when data cubes are having same/earlier dimensions.

e.g.: Sales



### 7) Range query

- Avg
- Sum
- Count
- Min
- Max

It applies a given aggregation operation over a set of selected cells.

For eg - find the total no. of sales for a period 2008-2010. Then the range query will return the value by aggregating the values of total sales from 2008-2010. It can be applied on dimensions as well.

numerical values.  
Range sum array [sum].

year	location	product	customers	Avg-sales	sold available
2008	Mumbai	Overs service	Ind Org	12 14	03
"	"	TV	Org	12	04.
"	"				01.