

Course Outcome

UNIT - I

→ Syllabus

→ Introduction :-

→ Data Warehouse And Mining

- A data warehouse is a subject oriented, Integrated, Non-volatile and time-variant collection of data in support of management's decision making process.

→ Textbooks

1. Jiawei Han and Micheline Kamber

"Data mining concepts and techniques"

Second edition Elsevier

Containing Def

- The four keywords, Subject oriented, Integrated, Time-variant and non-volatile distinguish Data warehouse from other data repository system such as Rational DB system, transactional processing system and file system

1. Subject Oriented

A data warehouse is organised around major subjects such as customers, suppliers, product and sales rather than concentrating on the day to day operations and transaction processing of an organization

2. Integrated

A data warehouse is usually constructed by integrating multiple heterogeneous sources such as relational database (DB), flat file and online transactional record.

3. Time-Variant

Data are stored to provide information from a historical perspective (the past 5-10 years).

4. Non-Volatile

A datawarehouse is always a physically separate store of data transform from the application data found in the operational environment. Due to this separation the data warehouse does not require transaction processing recovery and concurrency control mechanism.

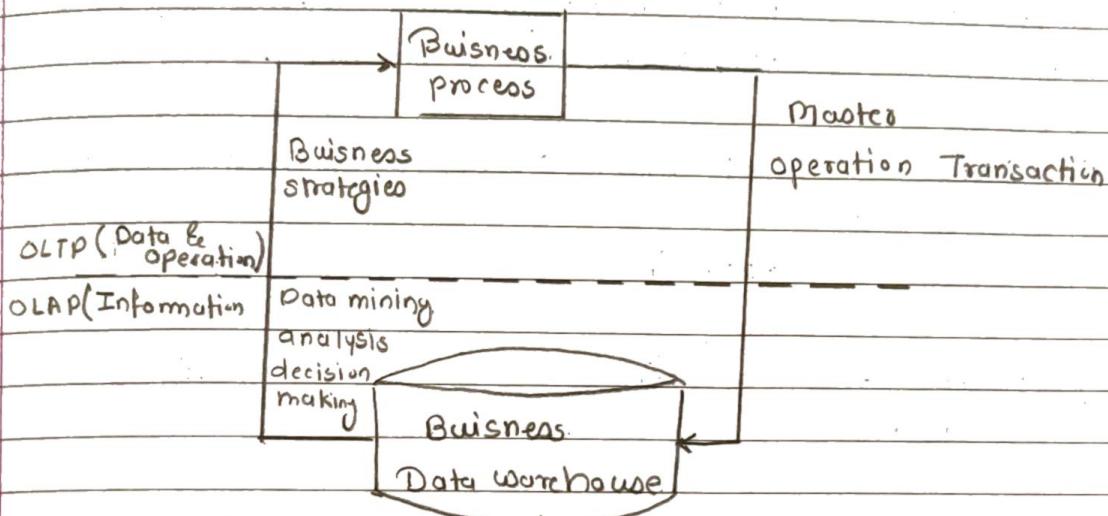
Data warehouse

It usually requires only two operation in data accessing i.e. initial loading of data and access of data.

* OLTP - Online Transaction processing

* OLAP - Online Analytic processing

25/10/29

OLAP VS OLTP

features	OLTP	OLAP
1) characteristics	operational processing	Information processing
2) Orientation.	Transcation	Analysis
3) User	clerk, DBA and any database professional	Knowledgeable worker. eg manager, executive & analyst
4) Function.	day to day transaction/operat operation	Long term information requirements, decision support
5) PB design	ER diagrams/ App oriented	Star/Snowflake, Subject oriented

6) Data	current data : guaranteed, upto date data	Historical data
7) Access	Read/write access	mostly read
8) Focus	data in	information.
9) Number of Records Access	Tens	millions
10) No. of Users	Thousands	Hundreds
11) DB size	100 MB to GB	100 GB to TB
12) priority	High performance, High availability	High flexibility, end users autonomy
13) Metric	Transaction throughput	query throughput, response time

3/01/20

Multidimensional Data Model.

- Data Warehouse and OLAP tools, are based on ~~the~~ multidimensional data model.
- The multi-dimensional data model is a method which is used for ordering data in the database.
- multidimensional data models views data in the form of data cube.
- Data Cube - A data cube allows data to be modeled, and viewed, in multiple dimensions. It is defined by dimensions & facts.
- Dimensions are the perspectives of the entities with respect to which an organisation want to keep records.

e.g. Take the example of data of a factory which sells products per quarter in location banglore.
 The data is represented in the table given below.

Time Quarter.	Location = Bangalore			
	Type of item	Bread	Sugar	Salt
Q1	350	389	35	50
Q2	260	528	50	90
Q3	483	256	20	60
Q4	436	396	15	40

2D - Factory Data.

In the above given presentation ; the factories sells for Bangalore are ,for the time dimension , which is organised into quarters and the dimension of items, which is sorted according to the kind of item which is sold.

The facts here are represented in Rs. (1000). Now, if we desire to view the data of the sells in a 3 dimensional table, then it is represented in the diagram given below.

Here the data of the sells is represented as a 2 dimensional table. Let us consider the data according to item , time and location.

Time	Location = "kolkata"			Location = "Delhi"			Location = "Mumbai"		
	Item	item	item	item	item	item	item	item	item
	Jam	Bread	Sugar	Jam	Bread	Sugar	Jam	Bread	Sugar
Q1	340	604	38	335	365	35	336	484	80
Q2	680	583	10	684	490	48	595	584	39
Q3	535	490	50	389	385	15	366	385	20

3D - data representation in 2D:

Mumbai	336	484	80
Delhi	335	365	35
kolkata			
Q1	340	604	38
Q2	680	583	10
Q3	535	490	50
	Jam	Bread	Sugar

→ Items (3D - representation)

The data can be represented in the form ³⁰ conceptually, which is shown in the image below.

5/02/23 Data Warehouse Schema :

1. Star schema.
2. Snowflake Schema.
3. Fact Constellation Schema.
4. star Schema.

time-Dimension Sales item
Table. Fact Table Dimension Table

time-key	time-key	item-key
day	item key	item-name
day-of-the-week	branch-key	brand
month	location-key	type
quarter	dollars-sold	Suppliers-key
year	units-sold	

Branch
Dimension Table

Location
dimension Table

branch-key	location-key
branch-name	street
branch-type	city
	Province-or-State
	country

Fig : star Schema

* Datawarehouse Schema

Schema is a logical description of the entire database. It includes the name and description of records of all record types including all associated data items and aggregates like entity-relationship data model which is commonly used in relational databases. In the design of where the databases schema consists of entities. a set of relationship b/w them.

The most popular data model for data warehouse is a multi-dimensional mode. Such a model can be exist in the form of the following types of schemas:

1. Star Schema
2. Snowflake schema.
3. Fact Constellation Schema

1. Star Schema

In Star Schema, data warehouse contains:

- 1) A large central table (fact table) containing the bulk of data, foreign keys with no redundancies.
- 2) A set of smaller attendant tables (dimension table), one for each dimension.

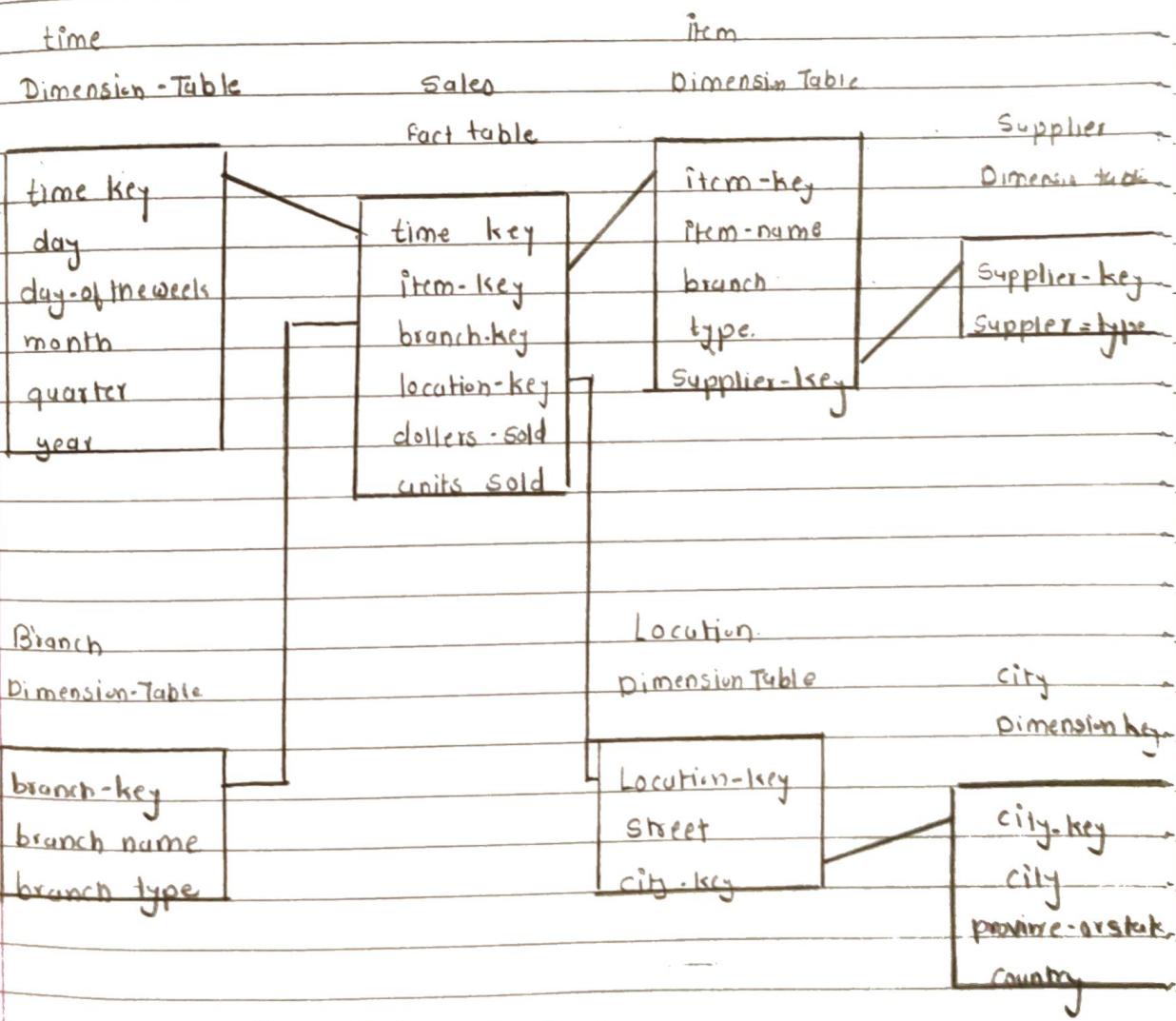
Consider the above diagram which shows the sales data of a company with respect to the 4 dimensions namely time, item, branch and location.

There is a fact table at the center. It contains the keys (foreign keys) to each of four dimensions. The

fact table also contains the attributes namely dollars sold and units sold.

In dimension table there can be data redundancy of data.

2- Snowflakes Schema



Snowflake Schema

In snowflakes schema... The main difference b/w the star and snowflakes schema is in the definition of dimension table... the single dimension table for item in the star schema is normalized in the snowflakes schema... resulting in new item and supplier tables. The normalisation splits up the data into additional table.

3. Fact - Constellation Table.

In fact - constellation table there can be multiple fact tables. This kind of schema can be viewed as a collection of stars and hence it is called as "Galaxy Schema".

Time

Dimension table

Item

Dimension table

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

Sales
Fact table.

Branch
Dimension Table

branch key
branch name
branch-type

time key
item key
branch key
location key
dollars sold
units sold

item key
item-name
brand
type
Supplier-key

Shipping
fact table

item key
time key
Shopper-key
from-location
to-location
dollars cost
units shipped

Location
Dimension Table

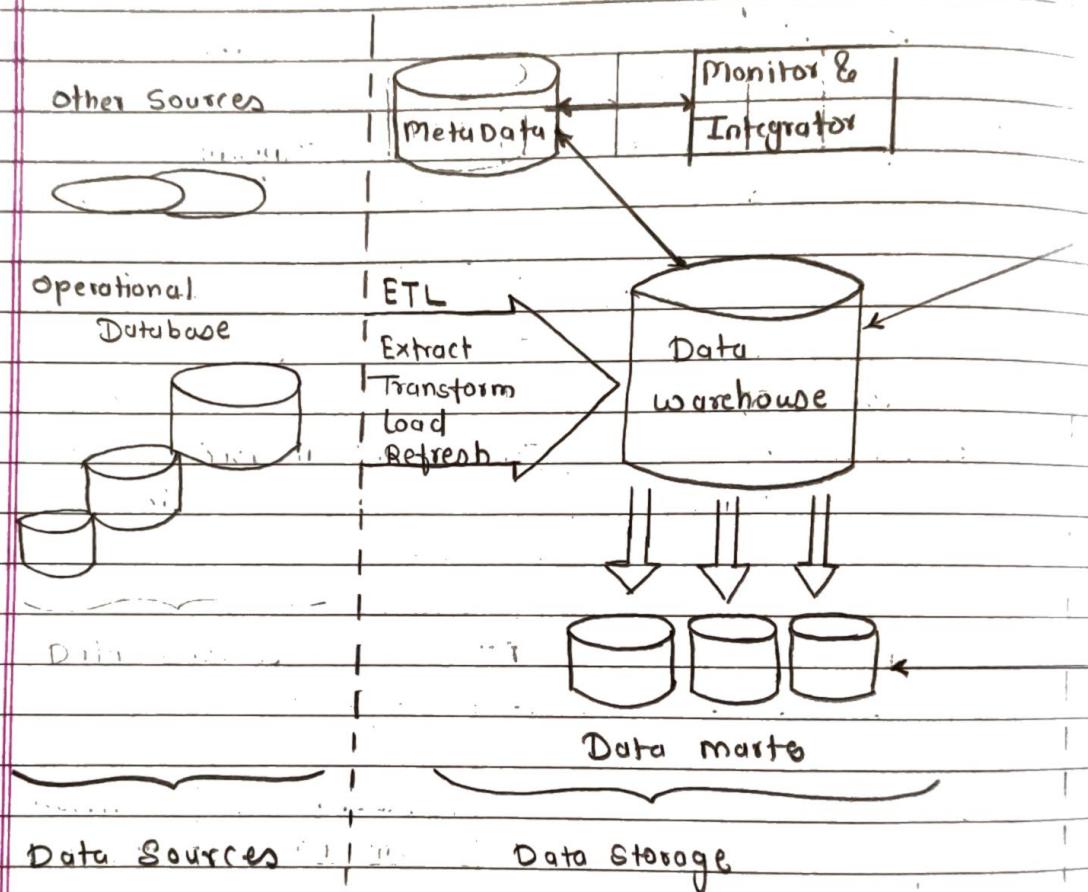
Location key
street
city - key
Province-or-state
country

Shipper
Dimension Table

Shopper - key
shipper key
location - key
shipper type

7/02/23

8-tier Data Warehouse Architecture

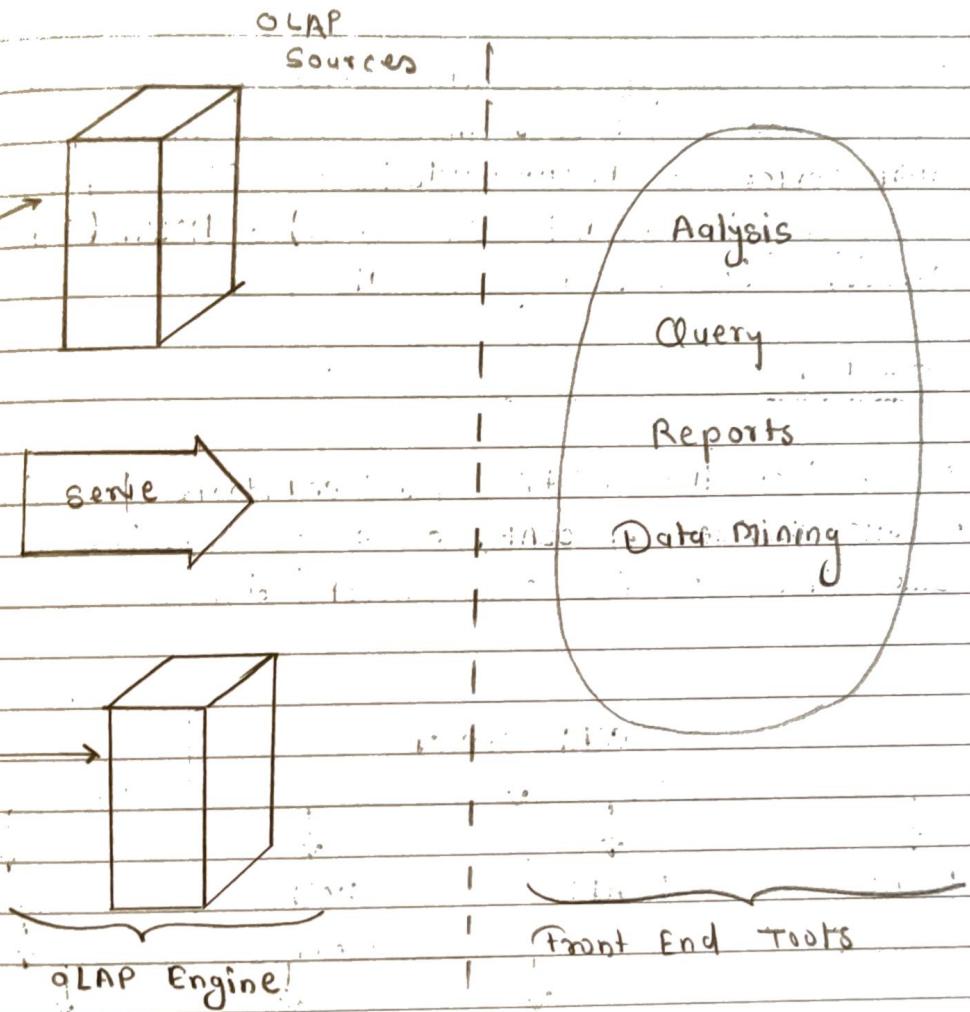


Data warehouse usually have 8 level (tier) Architecture. that includes

1) Bottom Tier

That consists of data warehouse server which is almost an RDBMS. It may include several specialised data marts and a metadata repository.

Data from operational databases and external sources (Such as user profile data provided by external consultant)



are extracted using application program interfaces called a gateway.

A gateway is provided by the underline DBMS and allows customer programs to generate SQL code to be executed at a server..

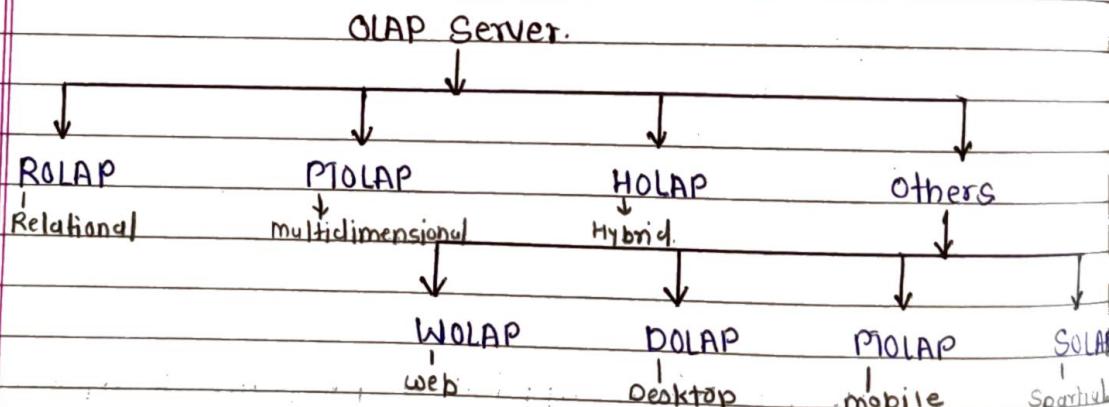
Examples of gateways contains ODBC (open database connectivity) and JDBC (Java database connectivity).

2) Middle tier.

It consists of an OLAP server for fast querying of the datawarehouse. There can be different types of OLAP server i.e. ROLAP (Relational OLAP server), MOLAP (multidimensional OLAP server), HOLAP (Hybrid OLAP server). (you have to write description)

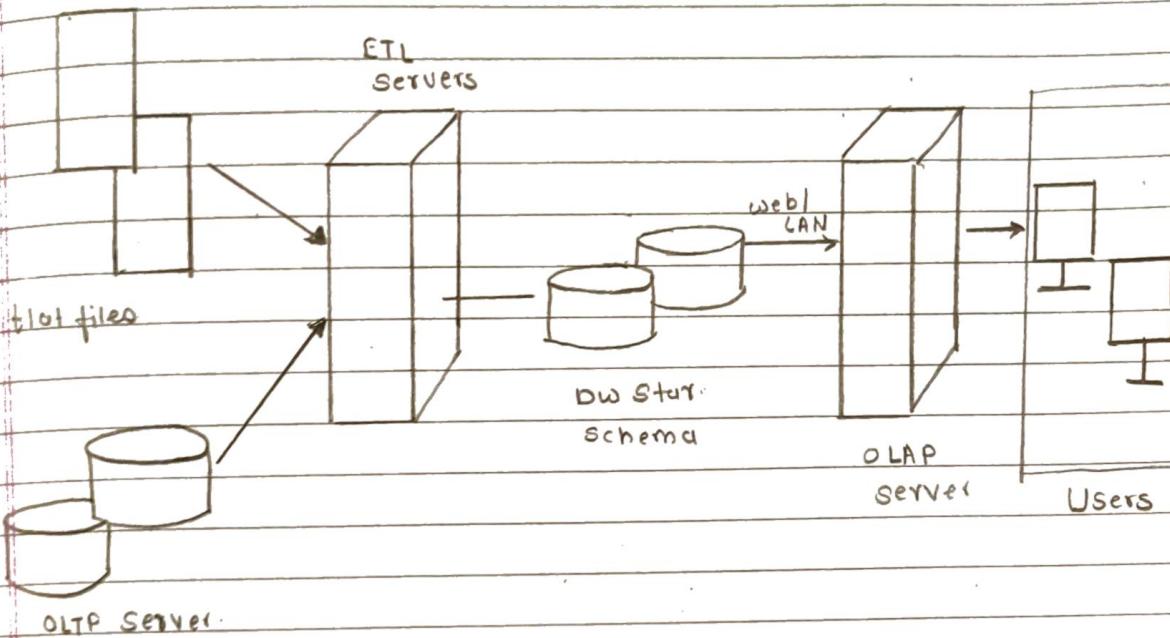
3) Top tier.

A Top tier that contains front end tools for displaying results provided by OLAP, as well as additional tools for data mining of the OLAP generated data.

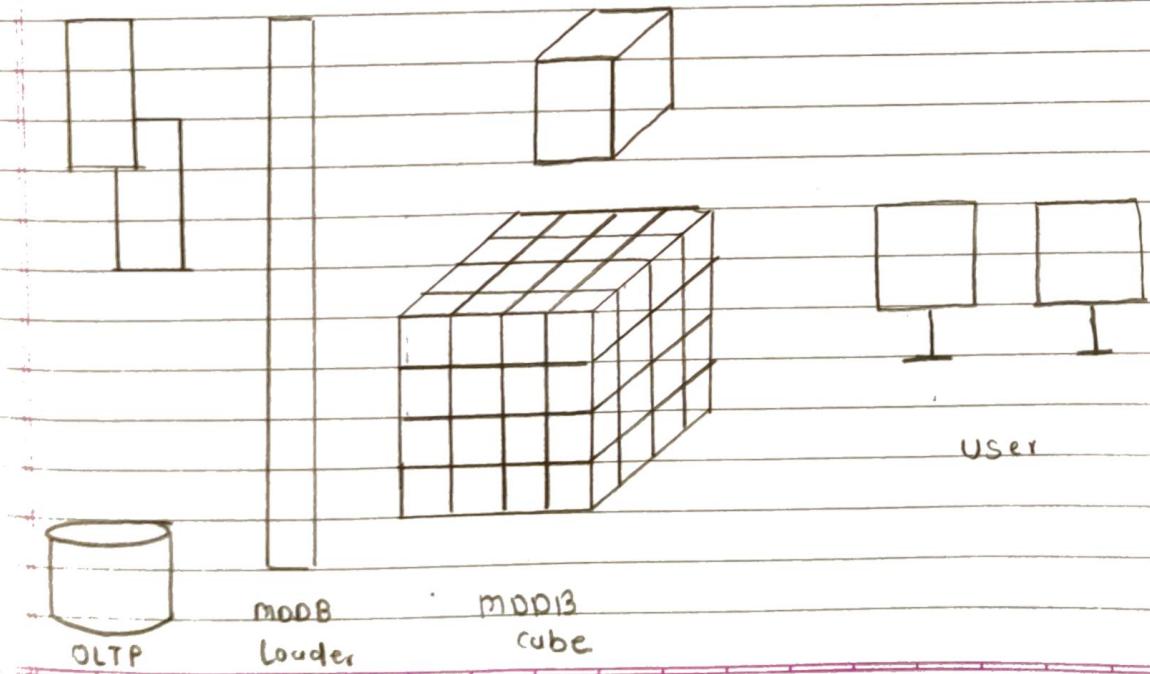


Types of OLAP server

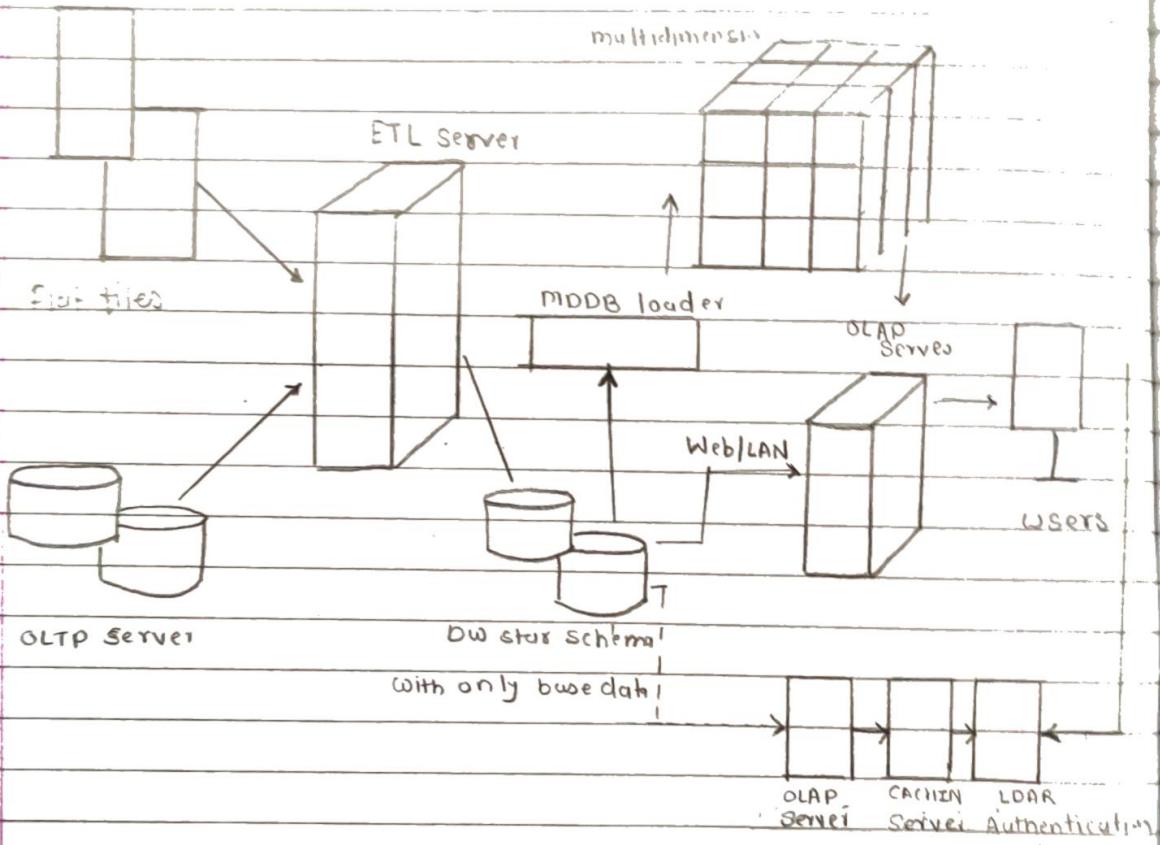
1) ROLAP Server



2) MOLAP Server



3) HOLAP Server



16/12/23

* Design & Construction of datawarehouse

In general warehouse design consists of following steps:

- 1) Choose a business process model.
- 2) choose the grain of the business process (fact table).
- 3) choose the dimension
- 4) choose the measures

Suppose that datawarehouse consist of 3 dimensions Time, doctor, patient and 2 measures i.e. count & charge where charge is the fee. That a doctor charges to the patient for a visit. Draw star schema diagram for the above datawarehouse.

Time
dimension table

Time-key
days
month
week
quarter

Fact Table	
Time-key	
Doctor key	
patient-key	
count	
charge	

Doctor dimension table

Doctor-key
Doctor-name
Phone-no.
Gender
Specialization

Patient dimension table

Patient-key
Patient-num
disease
illness
pharmacy

Suppose the datawarehouse consist of the 4 dimensions date, spectator, location and game and 2 majors count and charge where charge is the fair that a spectator pays for watching a game on a given date. Spectator may be student, adults or seniors with each spectator having its own charge. Draw the star schema.

date

dimension table

date-key
day
month
year
quarter

fact table

date-key
Spectator-key
location-key
game-key
count
charge

spectator

dimension Table

Spectator-key
Spectator-name
status
charge-date
phone-no
email-id

location

dimension table

location-key
city
state
country

game

dimension table

game-key
game-name
game-type
producer.name
actor

Unit 2Syllabus

Fundamental data mining

- 1) Data mining functionalities, classification of data mining system, data mining task primitives, major issues and challenge in data mining, data preprocessing / need for data preprocessing, data cleaning, integration, transformation, data reduction, data mining application areas.

* Data Mining

Data mining refers to extracting or mining knowledge from large amount of data.

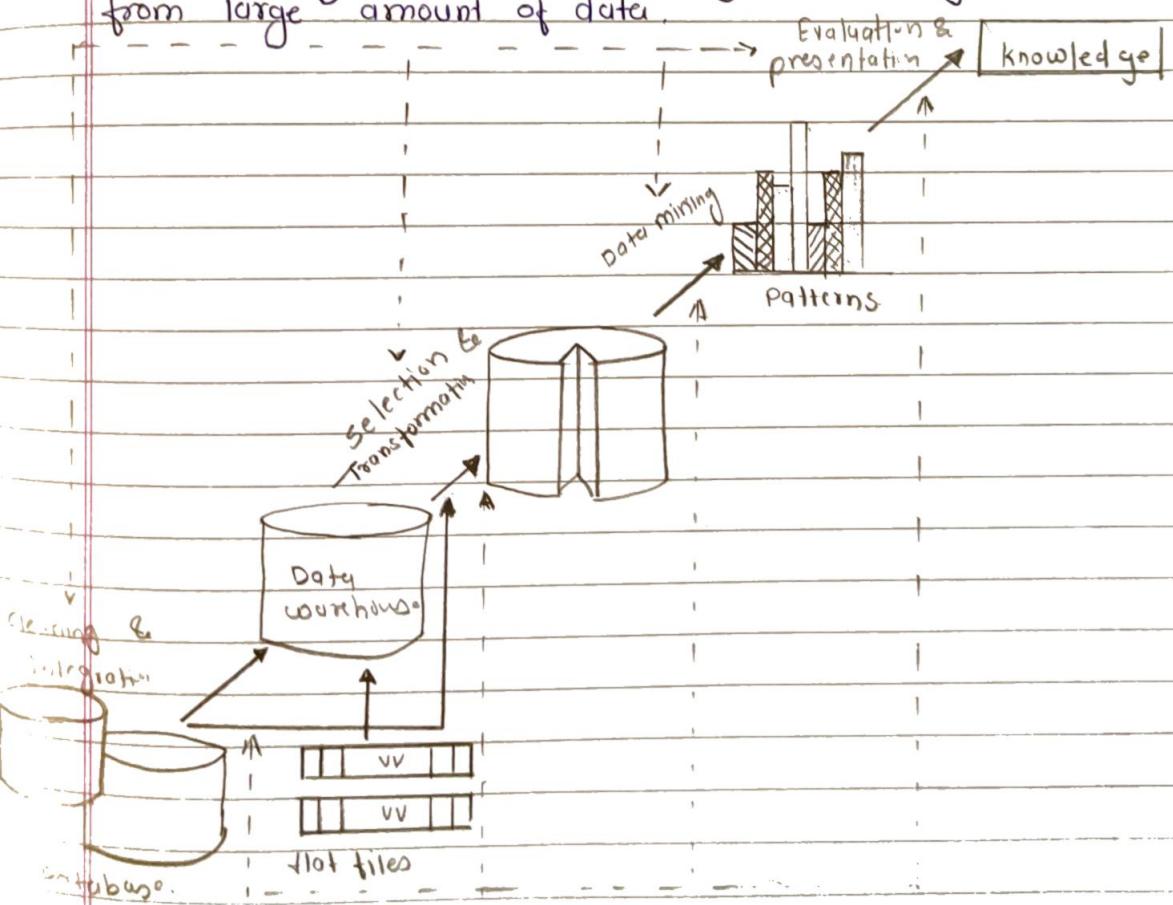


fig. Data mining as a step in the process of knowledge discovery

27/02/23

* Architecture of Data mining : major Components in Data mining

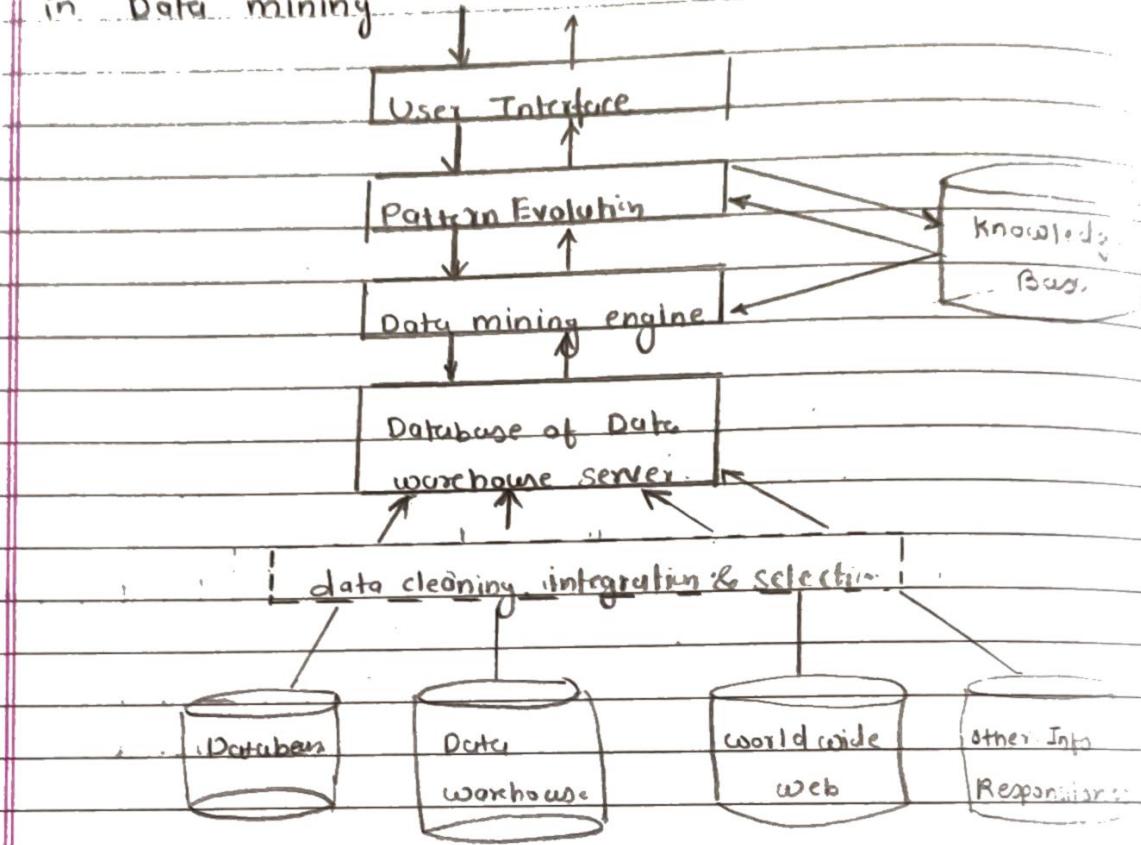


Fig Architecture of Typical Data mining System

* Data Mining Functionalities

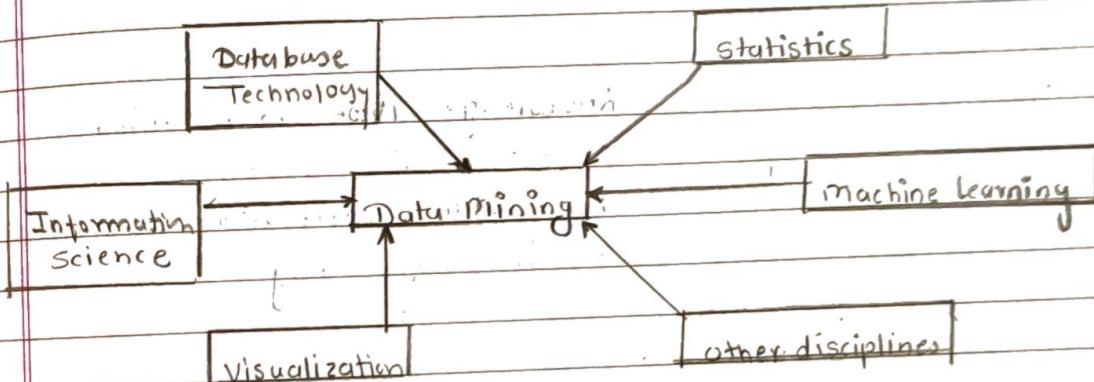
- 1) Concept / class Discription.
- 2) Mining Frequent Patterns , association , correlation
- 3) classification and prediction
- 4) Cluster Analysis

5) Outlier Analysis

6) Evolution Analysis

8/10/23
★

Classification of Data Mining System

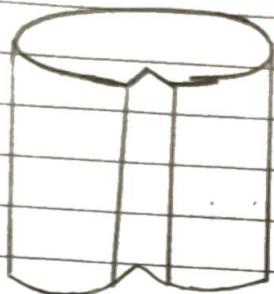


Data mining as confluence of multiple disciplines

- 1) classification according to kinds of database mined
- 2) classification according to the kinds of knowledge mines
- 3) classification according to the kinds of technique utilized
- 4) classification according to the applications adopted.

* Data Mining Task Primitives

1)



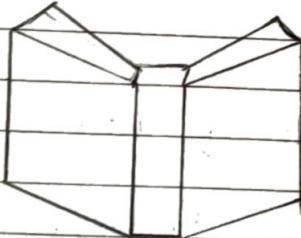
- 1) Task relevant Data
- 2) Database or Data warehouse name
- 3) Database Table or Data cubes
- 4) Condition for Data Selection
- 5) Relevant attributes or dimensions
- 6) Data Grouping criteria

2)



knowledge type to be mind
Characterization, Discrimination,
Association, correlation, classification,
prediction, clustering

3)



Background knowledge
concept hierarchies user
Beliefs about relationships
in the data

4)



Pattern interestingness measures
simplicity
certainty (e.g. confidence)

5)



Visualization of discussed
patterns

Rules, Tables, reports, charts, graphs,
decision trees & cubes etc
Drill down and roll up.