

→ Data mining architecture :-

The significant components of data mining systems includes :-

- Data ~~store~~ Sources.
- Data mining engine
- Data warehouse server
- Pattern evaluation
- gui
- Knowledge base

Data source:-

The actual source of data is database, data warehouse, world wide web(ww), text files and other document.

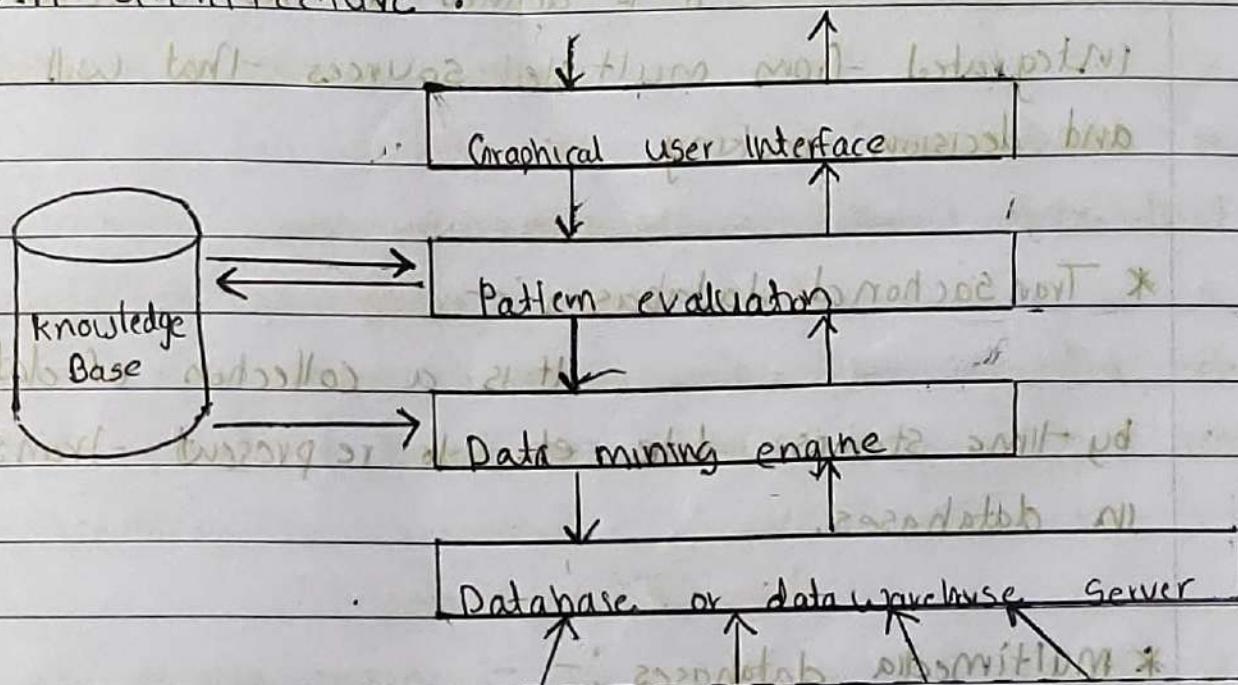
* Organization typically stores data in database or data warehouse.

Some of the data sources may be:-

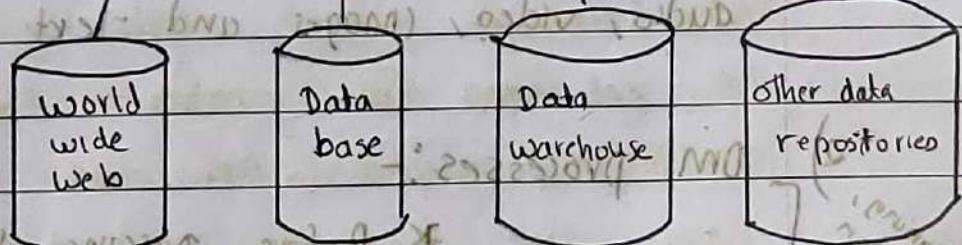
→ flat files:-

Data files in text form or binary form with a structure that can be easily extracted by data mining algorithms.

⇒ DM architecture :-



Data cleaning, Integration & Selection



⇒ Datamining components :-

1) Data Sources :-

* Relational databases :- A relational database is defined as

the collection of data organized in

tables with rows and columns.

Latin words such as column, most often used to

* Data Warehouse :-

It is defined as the collection of data integrated from multiple sources that will support queries and decision making.

* Transactional databases :-

It is a collection of data organized by time stamps, date, etc. to represent transaction in databases.

* Multimedia databases :-

Multimedia databases consists audio, video, images and text media.

DM processes :-

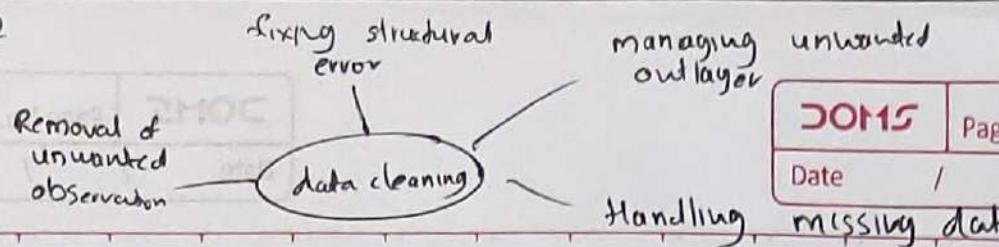
* Before passing the data to the database or data warehouse server, the data must be cleaned, integrated and selected.

* As the information comes from various sources and in different formats, it can't be used directly for the data mining procedure because the data may not be complete and accurate.

* So the first step requires to be cleaned and unified.

* Data integration in data mining refers to the process of combining data from multiple sources into single, unified

Data cleaning
Data integration



DOMS

Page No.

Date / /

view. after integrating the data we need to select the relevant data. that process is called data selection.

3) Database or Data Warehouse server :-
The database or data warehouse server consists of the original data that is already ready to be processed.
Hence the server is caused for retrieving the relevant data that is based on data mining as per user request.

4) Data mining engine :-
The data mining engine is a major component of any data mining system.
It contains several modules for operating data mining tasks, including association, classification, clustering, prediction, time-series analysis, etc.
It comprises instruments and software used to obtain insights and knowledge from data collected from various data sources and stored within the data warehouse.

5) Pattern evaluation module :-
The pattern evaluation module is primarily responsible for the measurement of investigation of the pattern by using a threshold value.

- * It collaborates with the data mining engine to focus the search on exciting patterns.
- * The pattern evaluation module might be coordinated with the mining module, depending on the implementation of both of the data mining techniques used.
- * for efficient data mining, it is abnormally suggested to push the evaluation of pattern stakes as much as possible into the mining procedure to confine the search to only fascinating patterns.

6)

Graphical user interface :-

- * The graphical user interface (gui) module communicates b/w the data mining system and the user.
- * This module helps the user to easily and efficiently use the system without knowing the complexity of the process.
- * This module cooperates with the data mining system when the user specifies a query or a task and displays the results.

7)

Knowledge Base :-

- * The knowledge base is helpful in the entire process of data mining.
- * It might be helpful to guide the search or to evaluate

The veracity of the resulting patterns in the knowledge base may even contain user views and data from user experiences that might be helpful in the data mining process.

The data mining engine may receive inputs from the knowledge base to make the results more accurate and reliable.

⇒ spatial database:-

store geographical information. Store data in the form of coordinate, topology, line, polygon etc.

⇒ KDD (Knowledge discovery in database)

KDD is a process involves the extraction of useful, previously unknown and potentially valuable information from large datasets.

It is a field of interest to researchers in various fields, including -

- Artificial intelligence
- Machine learning
- pattern recognition

- Database

- Knowledge acquisition for expert systems
- data visualization

- The KDD process is an iterative process and it

requires multiple iterations of the various steps to extract accurate knowledge from the data.

steps:-

1) Data cleaning :-

* Data cleaning is also known as data cleansing or data scrubbing, is the process of identifying and correcting (or removing) errors, inconsistencies and inaccuracies within a dataset.

* It can be done by filling missing value and smoothing ^(unwanted data) noisy data ~~(error)~~, analyzing and removing outliers and removing inconsistencies in the data.

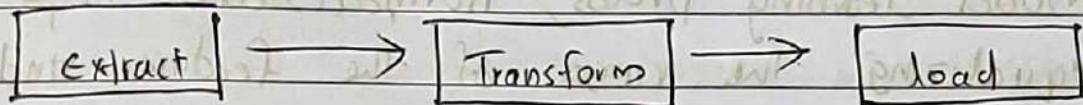
* The inconsistency can be recorded in various transactions during data entry or arising from integrating information from multiple database.

2) Data Integration :-

process of combining data from multiple sources to create a unified dataset.

* The heterogeneous data from multiple sources are combined into one common source.

* It is also known as ETL (Extract - Transform - Load) process.



iii) Data selection:-

It is defined as the process where data relevant to the analysis is decided and retrieved from the data collection.

iv) Data transformation:-

It is defined as the process of transforming data into appropriate form required by mining procedure.

The common method used for transformation are:-

- i) discretization
- ii) standardization
- iii) normalization.

→ normalization:-

• It is used to rescale the features to a standard range of values which is usually 0-1.

- ~~Discretization~~ → Normalization is usually used when different features have different range of values.

and some feature might contribute more to the model learning process, normalization helps in equalizing the range of the features and makes sure that the features contributes equally to the learning algorithm.

- 2. minimum and maximum value of features are used for scaling.

\Rightarrow Standardization:-

It is used to transform the data to standardized format.

- 3. Mean and standard deviation is used for scaling.

\Rightarrow Discretization:-

It is used to converting a huge number of data values into smaller ones so that the evaluation and management of data become easy.

\Rightarrow Data mining:-

It is defined as techniques that are applied to extract patterns potentially useful.

It transforms task relevant data into patterns, and decides purpose of model using classifiers or

characterization.

similar by some similarity exist between them.

v) pattern evaluation:-

It is the process of assessing the quality of discovered patterns.

- * This process is important in order to determine whether the patterns are useful and whether something can be trusted.

vi) knowledge representation:-

Is the presentation of knowledge to the user for visualization in terms of trees, tables, rules, graphs, charts, matrices, etc.

- * It is a place where visualization and knowledge representation techniques are used to present mined knowledge to users.

vii) KDD vs Data mining:-

KDD

- * refers to a process of identifying valid, novel, potentially useful and ultimately understandable patterns and relationships in data.

- * Technologies used: Data

* data mining is a step with it
End process.

data mining

- * Data mining refers to a process of extracting useful and valuable information or patterns from large data sets.

- * Technologies used:

cleaning, data integration, data selection, data transformation, data mining, pattern evaluation and knowledge representation and visualization.

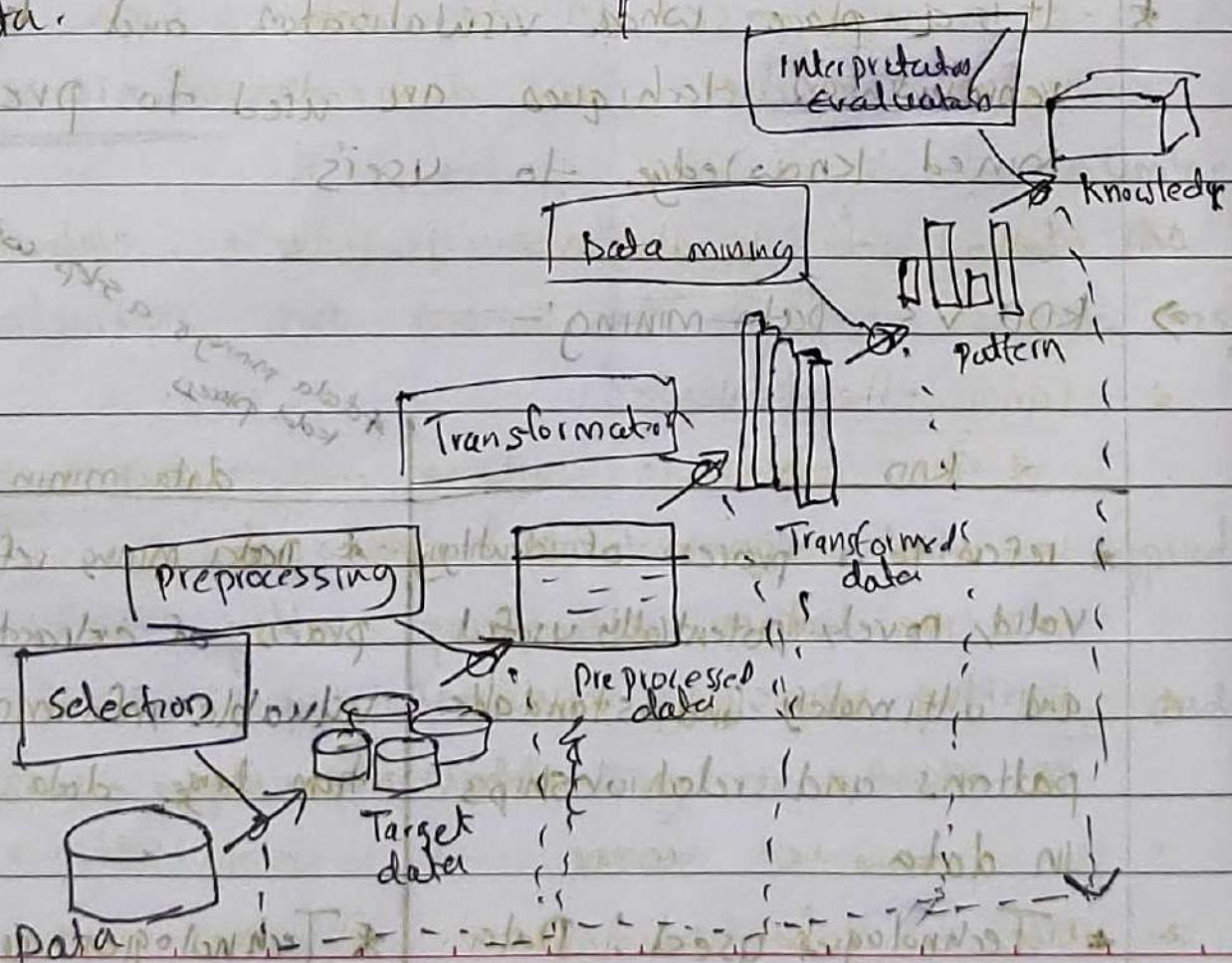
* Output:— structured information such as rules and models, that can be used to make decisions or predictions

* Focus is on the discovery of useful knowledge / rather than simply finding patterns in data.

Association rules, classification, clustering, regression, decision-tree, neural networks, and dimensionality reduction.

* Output:— patterns, associations or insights that can be used to improve decision-making or understanding.

* Data mining focus is on the discovery of pattern or relationship in data.



⇒ Data preprocessing:-

It is the process of converting raw data into an understandable format.

Steps:-

- 1) Data cleaning
- 2) Data integration
- 3) Data transformation
- 4) Data reduction (selection)

⇒ Data cleaning:-

This involves identifying and correcting errors or inconsistencies in the data, such as missing values, outliers, and duplicates.

⇒ Methods for handling missing data:-

* Ignore the tuples:-

→ This approach is suitable only when the dataset we have is quite large and multiple values are missing within a tuple.

* Fill the missing value:-

→ Fill the missing values manually or by the most probable value.

* Use a global constant to fill the missing value:-

- * In this replace all the missing attribute by a same constant.
 - * Use a measure of central tendency for the attribute to fill the missing value.
 - * In this replace all the missing attribute by central tendency values such as mean or median.
- ⇒ Methods for handling noisy data
- * Noisy data is a meaningless data that can't be interpreted by machines.
 - * It can be generated due to faulty data collection, data entry error etc.
 - * Steps - , Binning
 - Regression
 - clustering

Binning:

The data is divided into several segments of equal size. After that, the different methods are executed to complete the task.

Regression:

Data can be made smooth by fitting it to a regression function. The regression used may be,

linear (having one independent variable) or multiple (having multiple independent variables).

⇒ clustering (outlier analysis):-

clustering group the data in a cluster. Then the outliers are detected with the help of clustering.

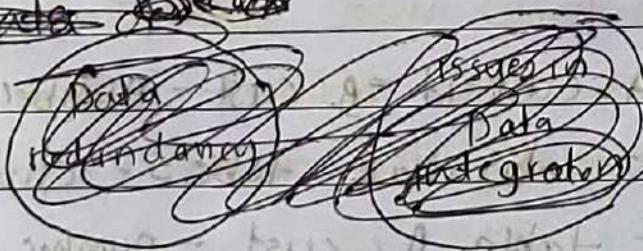
⇒ Data Integration:-

It is data preprocessing technique that combines data from multiple heterogeneous data sources into a coherent data store and provides a unified view of the data.

The process involves identifying and accessing the different data sources, mapping the data to a common format.

⇒ Issues in data integration:-

* Data redundancy



Data redundancy: Redundant data occurs when we

merge data from multiple databases.

* If the redundant data is not removed, incorrect results will be obtained during data analysis.

⇒ Duplicate data attributes:- (multiple occurrences of a year)

Identify the duplicate data attributes and remove those attributes.

⇒ Irrelevant attributes:-

Some attributes in the data are not important, and they are not considered while performing the data mining task.

⇒ Entity identification problem:-

Integration of data from multiple resources, some data resources match each other and they will become redundant if they are integrated.

* for example - A.cust-id = B.cust-number. Here A, B are two different database tables. cust-id is the attribute of table A, cust-number is the attribute of table B.

* Here cust-id and cust-number are attributes of different tables and there is no relationship

b/w these tables but the cust-id attribute and cust-number attribute are taking the same values.

- * It helps in detecting and resolving data value conflicts.

⇒ Data transformation:-

- This involves converting the data into a suitable format for analysis.
- Common techniques used in data transformation include. The different method involves:-
 - Normalization
 - Standardization
 - Discretization

⇒ Methods for data normalization :-

- Decimal Scaling
- min-max normalization
- z-Score Normalization
- (zero-mean normalization)

⇒ Data reduction:-

data reduction is a crucial step in the data mining process that involves reducing the size of the dataset while preserving the important information.

This is done to improve the efficiency of data analysis and to avoid over fitting of the model.

- Data cube aggregation
- Attribute subset selection
- Dimensionality reduction

⇒ Data cube aggregation :-

The technique is used to aggregate data in a simpler form.

- Data cube aggregation is a multi-dimensional aggregator that uses aggregation at various levels of a data cube to represent the original data set, thus achieving data reduction

⇒ Attribute Subset Selection :-

The large data set has many attributes, some of which are irrelevant to data mining.

- * The core attribute subset selection reduces the data volume by eliminating redundant and irrelevant attributes.
- * The attribute subset selection ensures that we get a good subset of original attributes even after eliminating the unwanted attributes.

⇒ Attribute selection

of observations. Whenever we find median value there are two intervals below median & above median.

$\frac{n}{2}$ th observation

$=$ n th observation

subset

set 20

subset

(char)

attribute

added to a reduced set.

* ~~stepwise~~ ^{step wise} ~~selection~~ ^{backward elimination} :-

- Here all the attributes are considered in the initial set of attributes
- In each iteration, one attribute is eliminated from the set of attributes whose p-value is higher than significance level.

* Combination of forward selection and backward elimination :-

- The stepwise forward selection and backward elimination are combined so as to select the relevant attributes most efficiently.

- This is the most common technique which is generally used for attribute selection.

⇒ Attributes subset selection methods :-

* Stepwise Forward Selection :-

- This procedure starts with an empty set of attributes as the minimal set.
- The most relevant attributes (having minimum p-value) are chosen and are added to the minimal set.
- In each iteration, one attribute is added to a reduced set.

~~step wise~~

→ backward elimination :-

- Here all the attributes are considered in the initial set of attributes.
- In each iteration, one attribute is eliminated from the set of attributes whose p-value is higher than significance level.

* Combination of forward selection and backward elimination :-

- The stepwise forward selection and backward elimination are combined so as to select the relevant attributes most efficiently.
- This is the most common technique which is generally used for attribute selection.

c) Dimensionality reduction :-

The number of input features, variables, or columns present in a given dataset is known as dimensionality, and the process to reduce these features is called dimensionality reduction.

- It is a way of converting the higher dimensions dataset into lesser dimensions dataset ensuring that it provides similar information.

⇒ Measures of central tendency :-

* mean

* median

* mode

* mid range

* range

⇒ mean :-

mean of data = sum of observations / no. of observations

Ex:- If there are 5 numbers 1, 2, 3, 4, 5 then mean is $\frac{1+2+3+4+5}{5} = 3$

Ex:- If there are 5 numbers 1, 2, 3, 4, 5 then median is 3.

⇒ median :-

means

median → to identify the middle value

from a set of observations. Whenever we need to calculate median value there are two conditions,

when n is odd,

$$\text{median} = (n+1)/2 \text{ in observation}$$

$$\underline{\text{eq}} - \underline{q_{11} = \frac{q+1}{2} = \frac{10}{2} = 5^{\text{th}} \text{ observation}}$$

~~when n is even~~

$$\text{definition of median} = \frac{[(n/2)^{\text{th}} \text{ observation} + (n/2+1)^{\text{th}} \text{ observation}]}{2}$$

⇒ Mode :-

most frequently occurring value in the data

$$\text{eg. } \text{and } \begin{pmatrix} 3 & 3 & 3 \end{pmatrix} = 7 \quad 7 \quad 9 \quad 14 \quad 13 \quad 16$$

here mode is 3/4 nondirectional

$$\sum_{\Omega \in \mathcal{F}} (\bar{x} - x) = 0$$

Q Find the mean of first 10 odd integers.

$$\text{odd} = 1, 3, 5, 7, 9, 11, 13, 15, 17, \frac{19}{\cancel{19}}, 19$$

birds moon ~~even~~^{31 NOV} sum of observation

No. of observations Answers?

$$= \frac{160}{17} > 1011$$

⇒ mid-range :- ~~another measure to determine mid~~

~~out~~ ~~sum~~ ~~it is defined as the sum of maximum and minimum value divided by 2.~~

$$\text{mid range} = \frac{\text{max} + \text{min}}{2}$$

⇒ range :-

~~range~~ ~~(1)~~ ~~it is the difference b/w the maximum and minimum values of observation in the data.~~

$$\text{range} = \text{maximum value in data} - \text{min value in data}$$

⇒ Variance :- ~~another measure to determine~~

~~variance of data is given by measuring the distance of observed values from the mean of distribution.~~

$$\sigma^2 = \sum_{i=0}^n (x_i - \bar{x})^2$$

~~explanation b/w of trait to norm with trait~~

⇒ standard deviation :- ~~another measure to determine~~

~~Square root of variance~~ ~~is called standard deviation~~

$$\sigma = \sqrt{\frac{\sum_{i=0}^n (x_i - \bar{x})^2}{n}}$$

Q. What is the median of following data set?

32, 6, 21, 10, 8, 11, 12, 36, 17, 16, 15, 18, 40, 24, 21, 23, 24, 24, 29, 16, 32, 31, 10, 30, 35, 32, 18, 39, 12, 20

Find the mid-range and range of the data set.

Step 1: First to ascending order.

6, 8, 10, 10, 11, 12, 12, 15, 16, 16, 17, 18, 18, 20, 21, 23, 24, 24, 24, 29, 30, 31, 32, 32, 32, 35, 36

~~Step 2: Write the sorted values~~ 39 / 40

No. of values in data set $\Rightarrow n = 30$

x of individual bracket with best

$$\text{Median} = \frac{\cancel{20}/\cancel{2}}{2}, \frac{(n/2)^{\text{th}} + ((n/2)+1)^{\text{th}}}{2}$$

$$\frac{21+21}{2} = 21$$

$$\text{mid range} = \frac{\max + \min}{2} = \frac{40 + 6}{2} = \frac{46}{2} = 23$$

$$\text{range} = \max - \min = 40 - 6 = 34$$

$$\frac{1+(1/n) + n(1/n)}{2} = \text{mid range}$$

$$\frac{20}{5} =$$

Module-2Database:-

It is a collection of data and stores data in tables. It deals with operational or transactional data. It can store MB and GB of data. A student db can store details of students and retrieve information based on user query. Used for OLTP (Online Transaction process).

Data warehouse

Store huge amounts of data. Data collected from multiple heterogeneous sources like files, DBMS etc. It is mainly used to store historical data. Can store TB's of data. eg - how the placement of cs students has improved over the last 10 years. It is used for OLAP (online analytical processing). A data warehouse is a system that stores data from a company's operational databases as well as external sources.

- Data warehouse platforms are different from operational databases because they store historical information, making it easier for business leaders to analyze data over a specific period of time.

Module - 2Database:-

It is a collection of data and stores data in tables. It deals with operational or transactional data. It can store MB and GB of data. A student db can store details of students and retrieve information based on user query. used for OLTP (online Transaction process).

Data warehouse:-

Store huge amounts of data. Data collected from multiple heterogeneous sources like files, DBMS etc. It is mainly used to store historical data. Can store TB's of data eg - how the placement of CS students has improved over the last 10 years. It is used for OLAP (online analytical processing). A data warehouse is a system that stores data from a company's operational databases as well as external sources.

- Data warehouse platforms are different from operational databases because they store historical information, making it easier for business leaders to analyze data over a specific period of time.

father of data warehouse

- According to bill Inmon, "A warehouse is a subject oriented, integrated, time, variant and non-volatile collection of data in support of management's decision making process."

⇒ Characteristics of data warehouse :-

- Subject oriented
- Integrated
- non volatile
- Time variant.

(volatile means data lose during power loss)

⇒ Subject oriented:-

DW provide a ~~concise~~ view around a particular subject such as customer, product or sales, instead of the global organization's ongoing operations.

- DW is always a subject-oriented one, as it provides information about a specific theme.

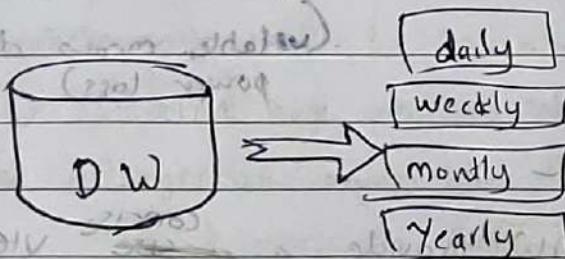
⇒ Integrated:-

DW integrates data from various heterogeneous data sources like ~~DBMS~~, flat files, and online transaction records, and combines it in a relational database.

- It requires performing data cleaning and integration during data warehousing to ensure consistency among different ~~dat~~ different data sources.
- It must be consistent, readable and coded.

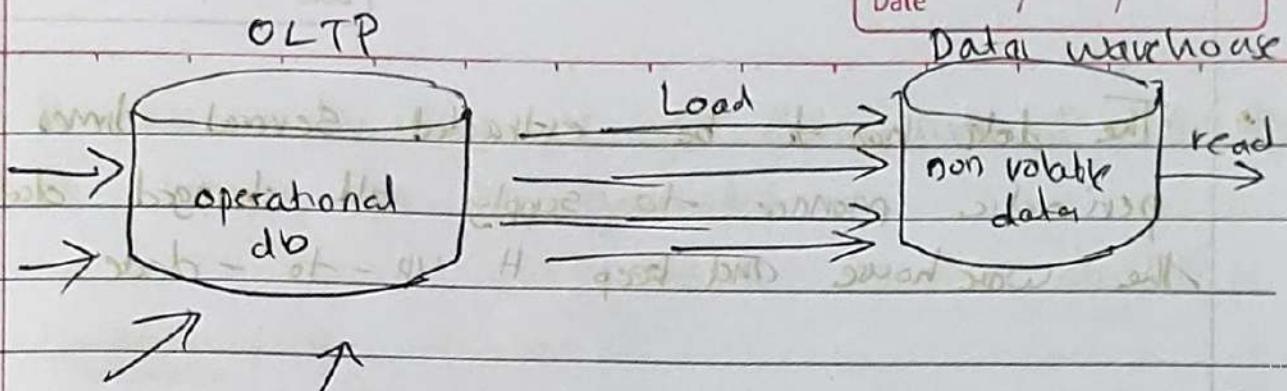
⇒ Time variant:

- Historical information is kept in a data warehouse.
- For example, one can retrieve files from 3 months, 6 months, 12 months or even previous data from a data warehouse.



⇒ Non volatile:

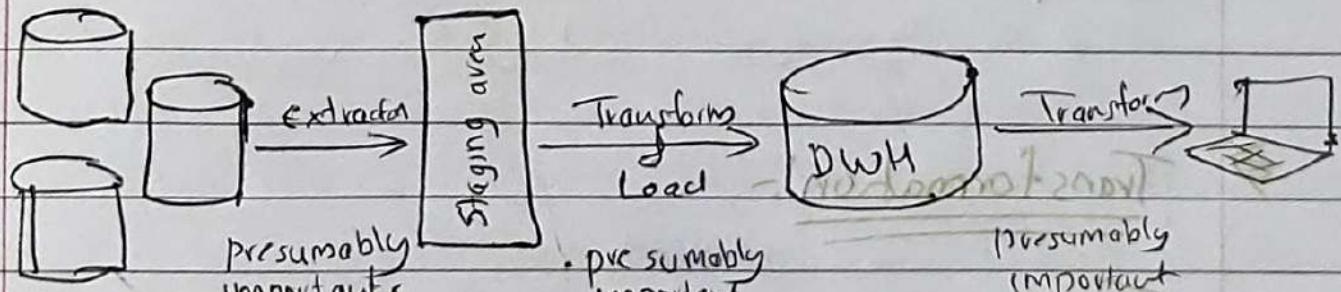
- The data residing in the data ~~warehouse~~^{Warehouse} is permanent.
- It ensures that when new data is added, data is not erased or removed.
- A data warehouse is kept separate from the operational db and thus the data warehouse does not represent regular changes in the operational database.



~~select / insert / delete / update~~

⇒ ETL process (extract, transform and load) :-

- The mechanism of extracting information from source systems and bringing it into the data warehouse is commonly called ETL.



⇒ Extraction:-

Extraction is the operation of extracting information from a source system for further use in a data warehouse environment.

- This is the first stage of the ETL process.

- The data has to be extracted several times in a periodic manner to supply all changed data to the warehouse and keep it up-to-date.

Cleansing:-

- The cleaning stage is crucial in a data warehouse because it is supposed to improve data quality.
- The primary data cleansing features found in ETL tools are rectification and homogenization.
- They use specific dictionaries to rectify typing mistakes and to recognize synonyms and defines appropriate associations b/w values.

Transformation:-

- Transformation converts records from its operational source format into a particular data warehouse format.
i.e. data extracted from source server is raw and not usable in its original form, therefore it needs to be cleansed, mapped and transformed.

- The data has to be extracted several times in a periodic manner to supply all changed data to the warehouse and keep it up-to-date.

⇒ Cleansing:-

- The cleaning stage is crucial in a data warehouse technique because it is supposed to improve data quality.
- The primary data cleansing features found in ETL tools are rectification and homogenization.
- They use specific dictionaries to rectify typing mistakes and to recognize synonyms and defines appropriate associations b/w values.

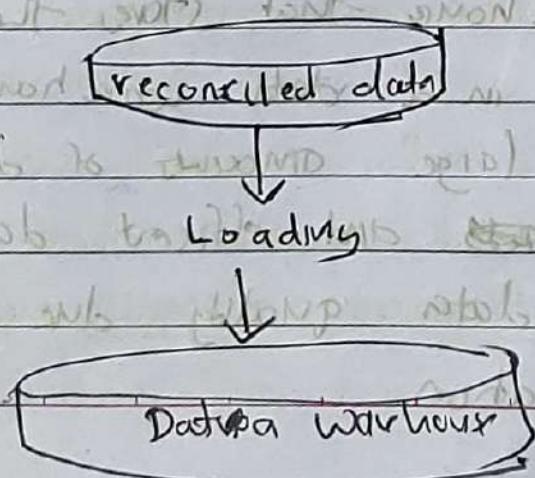
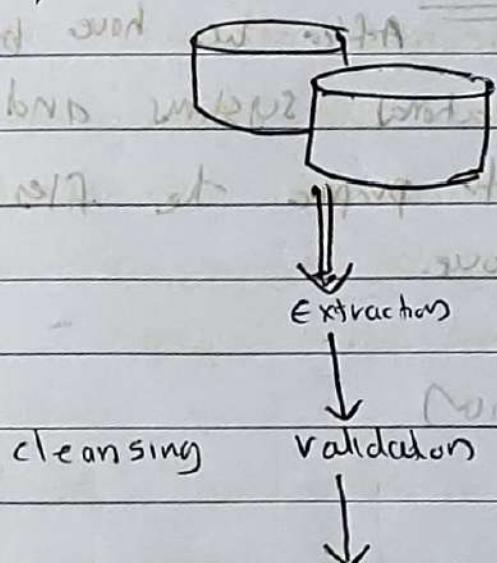
⇒ Transformation:-

- Transformation converts records from its operational source format into a particular data warehouse format.
i.e. data extracted from source server is raw and not usable in its original form, therefore it needs to be cleansed, mapped and transformed.

It adds value and changes data such that insightful BI reports can be generated.

⇒ Loading:-

- Loading data into the ~~target~~ target data warehouse is the last step of the ETL process.
- Huge volume of data needs to be loaded in a relatively short period.
 - Hence, load process should be optimized for performance.



⇒ Components or building blocks of data warehouse system

i) ⇒ Data source components :-

- i) Internal data
- ii) Archived data
- iii) External data

ii) ⇒ Data staging components :-

After we have been extracted data from various operational systems and external sources, we have to prepare the files for storing in the data warehouse.

i) Data extractors

ii) Data transformation

iii) Data Loading

iii) ⇒ Data storage component :-

Data storage is the component of the data warehouse that stores the data. Advantages of data storage in a data warehouse include the ability to store large amounts of data in a single location, fast and efficient data retrieval, and improved data quality due to data cleaning and standardization.

→ Data warehouse architecture :-

Data warehouse architecture defines the arrangement of the data in different databases. The different architectures include

- Single tier architecture
- Two tier architecture
- Three tier architecture.

→ Single tier architecture :-

An operator system is a method used in data warehousing to process the day-to-day transactions of an organization.

- A flat file system is a system of files in which transactional data is stored, and every file in the system must have a different name.
- Meta data summarizes necessary information about data which can be used to access data more easily.
- end-user access tools provides information to the business managers for strategic decision-making.
- The various end user tools includes:-
 - * reporting and querying tools
 - * Application development tools.
 - * Executive information system tools.
 - * online analytical processing tools.
 - * Data mining tools.

This image shows a close-up of a page from a spiral-bound notebook. The page is ruled with horizontal lines. Handwritten text is present, but it is very blurry and difficult to decipher. The handwriting appears to be in cursive script.

- A single tier - architecture helps to minimize the amount of data stored to reach the goal, i.e., it removes data redundancies.

⇒ Two tier architecture :-

 Application layer ←

 Source layer :- ←

A data warehouse system uses a heterogeneous source of data

- That data is stored initially to corporate relational databases or legacy databases, or it may come from an information system outside the corporate walls.

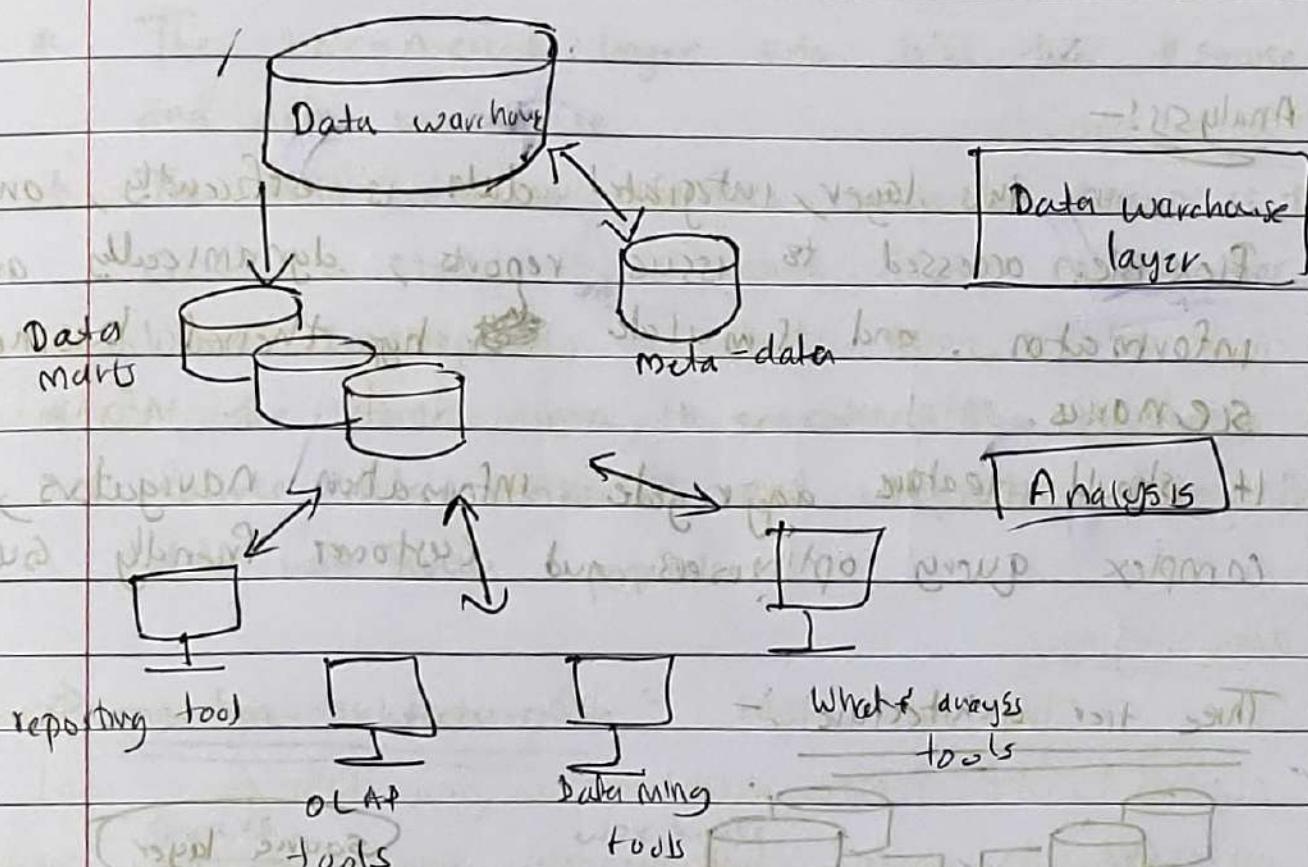
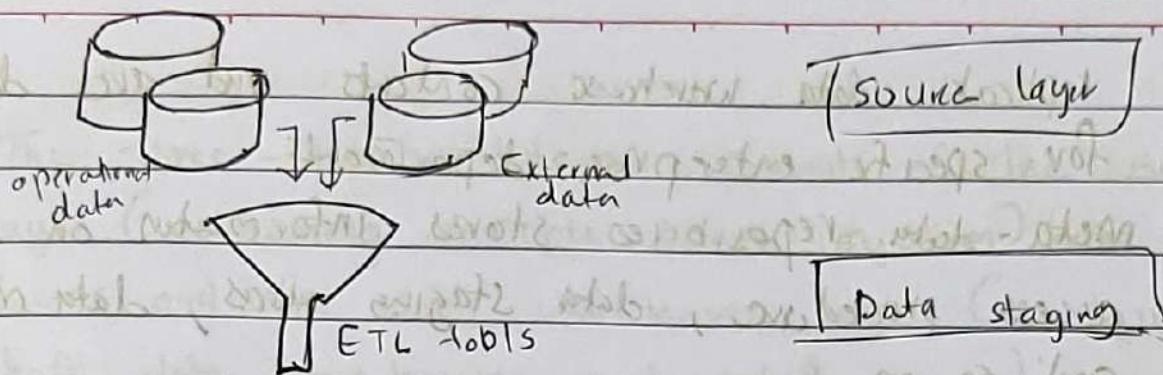
⇒ Data Staging:-

The data stored to the source should be extracted,

cleansed to remove inconsistencies and fills gaps, and

integrated to merge heterogeneous sources into one standard schema.

- Extraction, Transformation and loading Tools (ETL) can combine heterogeneous schemas, extract, transform, cleanse, validate, filter and load source data into a data warehouse.



⇒ Data warehouse layer:-

- Information is saved to centralized individual repository i.e., the data warehouse.
- The data warehouse can be directly accessed, but it can also be used as a source for creating data marts, which are which, partially ~~not~~.

Replicate data warehouse contents and are designed for specific enterprise departments.

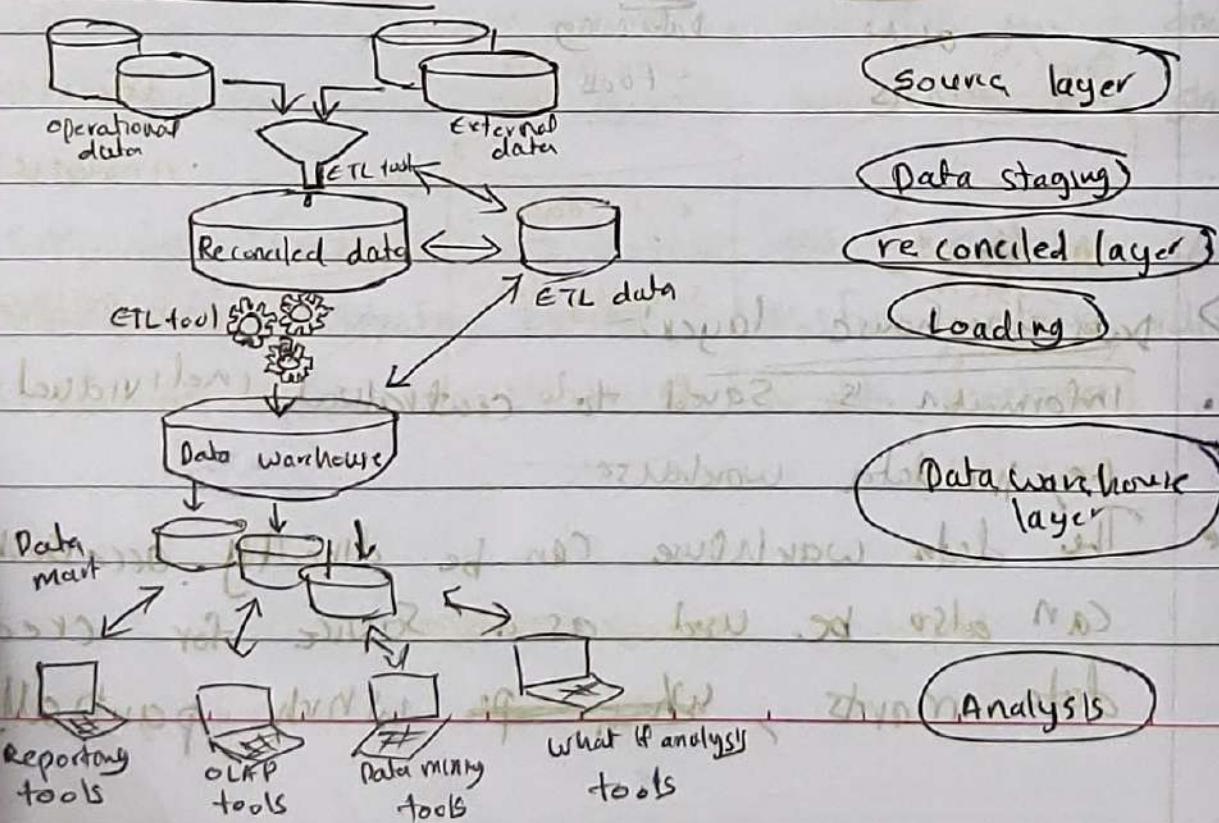
- Meta-data repositories stores information on sources, access procedure, data staging users, data mart schema and so on.

Analysis :-

In this layer, integrated data is efficiently, and flexibly accessed to issue reports, dynamically analyse information, and simulate hypothetical business scenarios.

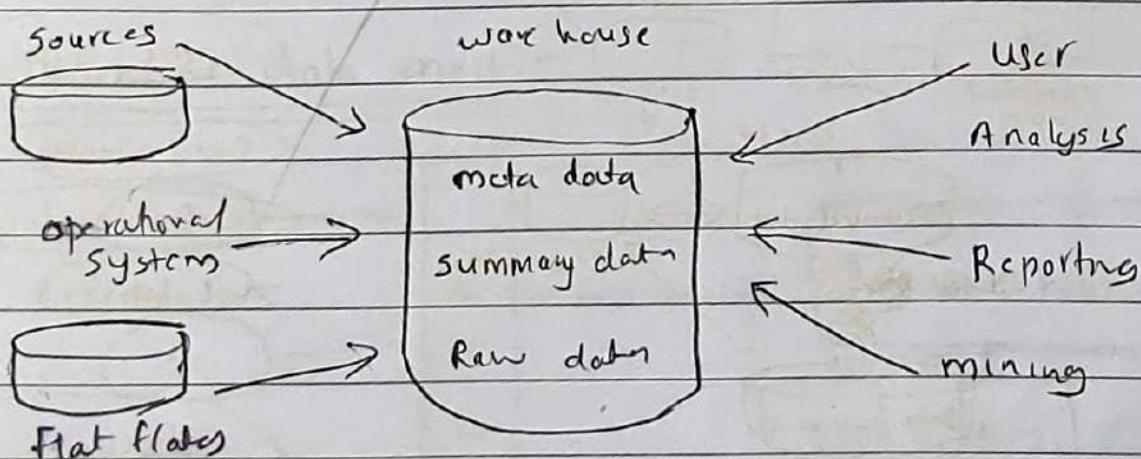
- It should feature aggregate information navigators, complex query optimizers and customer friendly GUI's.

Three tier architecture :-



- * The three-tier architecture consists of the source layer (containing multiple source systems), the reconciled layer and the data warehouse layer (containing both data warehouses and database marts):
- * The reconciled layer sits b/w the source layer and data warehouse.
- * The main advantage of the reconciled layer is that it creates a standard reference data model for a whole enterprise.
- * At the same time, it separates the problems of source data extraction and integration from those of data warehouse population

⇒ one tier architecture :-



Input → storage → o/p