

Data Mining

Unit-1

Data Mining :-

Data Mining is defined as procedure of extracting information from huge sets of data.

(or)

It is also defined as mining knowledge from data.

Data :-

Data is defined as facts (or) figures (or) information that's stored in (or) used by computer.

Types of Data :-

There are three types of data. They are :

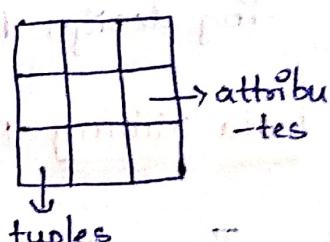
1. Database
2. Data ware House
3. Transactional database

1. Database Data (RDBMS) :-

It is set of tables, which contains rows and columns.

Row Represents the tuples

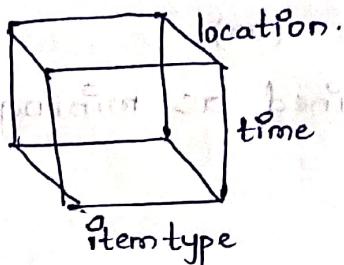
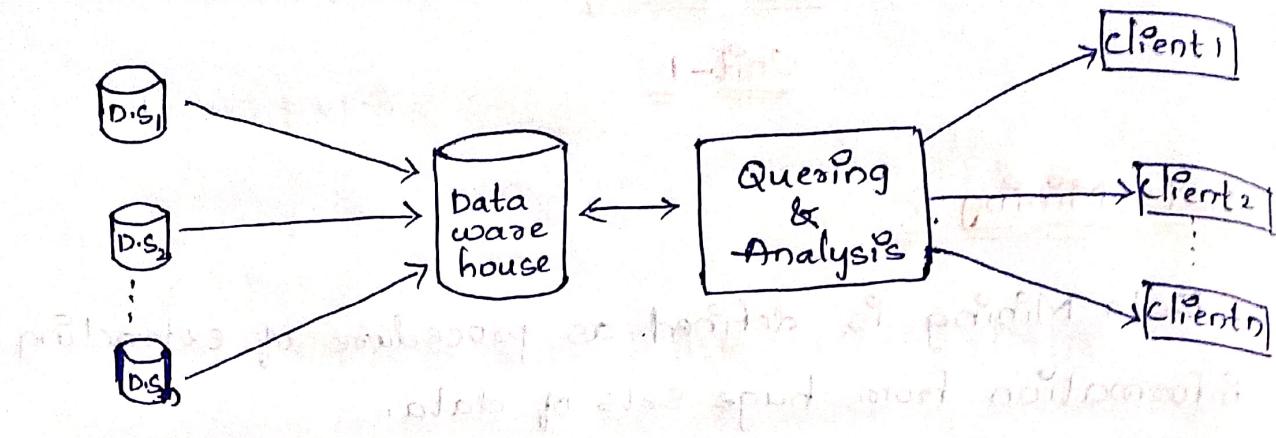
Column Represents the attributes.



2. Data ware House :-

Collection of data integrated from different sources with querying and decision making on data.

In data ware house, data is stored in dimensional structure (data cube) where each dimension is each attribute.



3. Transactional database: ~~part of transactional data~~

Each record is called as transaction

(sales, flight booking, etc)

Transaction has transaction ID, list of other items making transactions.

From transaction DB, we can mine frequent patterns

Other types of data:

Sequence data, spatial data, data streams, engineering design data, web data, etc.

Data Mining Functionalities :-

There are 5 functionalities. They are:

1. Concept / class descriptions

2. Mining frequent patterns, association & correlation

3. Classification & Regression for predictive analysis.

4. Cluster Analysis.

5. Outlier Analysis.

① Class Description / Concept:

Data is associated with class/concept.

Descriptions can be done in two ways.

a) Data Characterisation

b) Data Discrimination

a) Data Characterisation:

It refers to the summary of the class/concept.
O/p → General Overview.

b) Data Discrimination:

It compares the common features of class/concept.
O/p → barcharts, curves, etc.

② Mining frequent patterns, associations & correlation:

Frequent Patterns:

Things which are found most commonly in data.

Frequent Itemsets

Frequent Subsequence

Frequent Substructure.

Association Analysis:

It is a way of identifying the relation between various items.

Ex: used to determine sales of items that are frequently purchased together.

Correlation Analysis:

It is Mathematical technique.

It shows how strongly pair of attributes are related together.

Ex: Tall people tend to have more weight.

③ Classification and Regression for predictive analysis

classification :-

Process of finding a model that distinguishes data items.

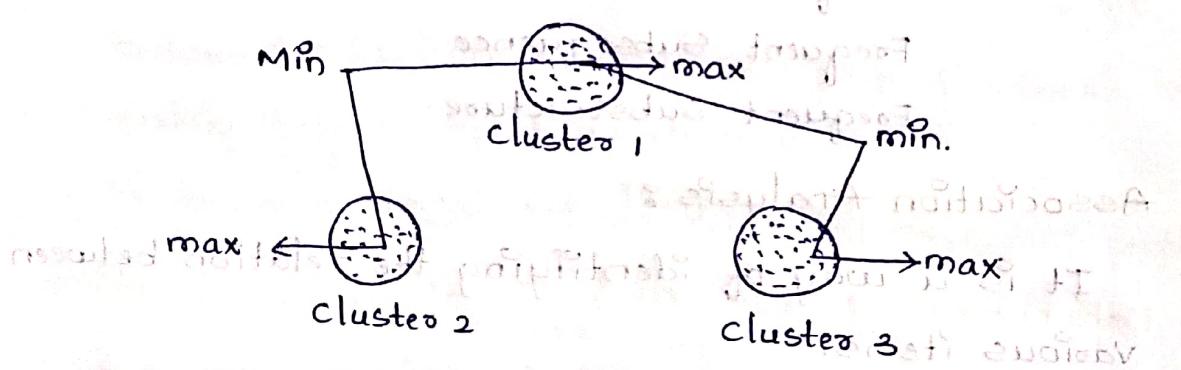
Decision tree is used for classification.

Regression :-

It is statistical methodology that is used for numeric prediction of missing data.

④ Cluster Analysis :-

The data items are distributed based on the principles of maximising the intracluster similarity and minimising the intercluster similarity.



⑤ Outlier Analysis :-

Among the data items in a DB, there may be some items which do not follow the behaviour of general data.

Those data items are called as outliers (noise/exceptions).

Interestingness of Patterns

In a data mining system, every day millions of data patterns are generated.

Among all these patterns generated, how many are really interesting?

Actually a small fraction of patterns generated would be of interest to any given user.

This raises three questions:

1. What makes patterns interesting?

A pattern is interesting if it is

- Easily understood by humans

- Validated new/test data.

- It is potentially useful.

2. Can data mining systems generate all of the interesting patterns?

- It refers to completeness of data mining system.

In reality it is not possible for a data mining system to generate all interesting patterns.

3. Can data mining systems generates only interesting patterns?

- It refers to optimization of a data mining system.

- It generates only interesting patterns, because it becomes easy and efficient for the user.

- Time is saved.

Classification of data mining Systems

Classification:

Process of finding a model that distinguishes the data items.

Decision trees are used for classification.

- Data Mining is everywhere and anywhere. It is used by many users.
- Data mining systems are classified based on several criteria.

1. Classification based on mined databases:

It is based on database that has been mined

- Relational
- Transactional
- Object relational

2. Classification based on type of knowledge mined:

- characterization
- discrimination
- Association and correlation analysis
- classification
- Prediction

3. Classification based on Outlier Analysis

4. Classification based on evolution Analysis

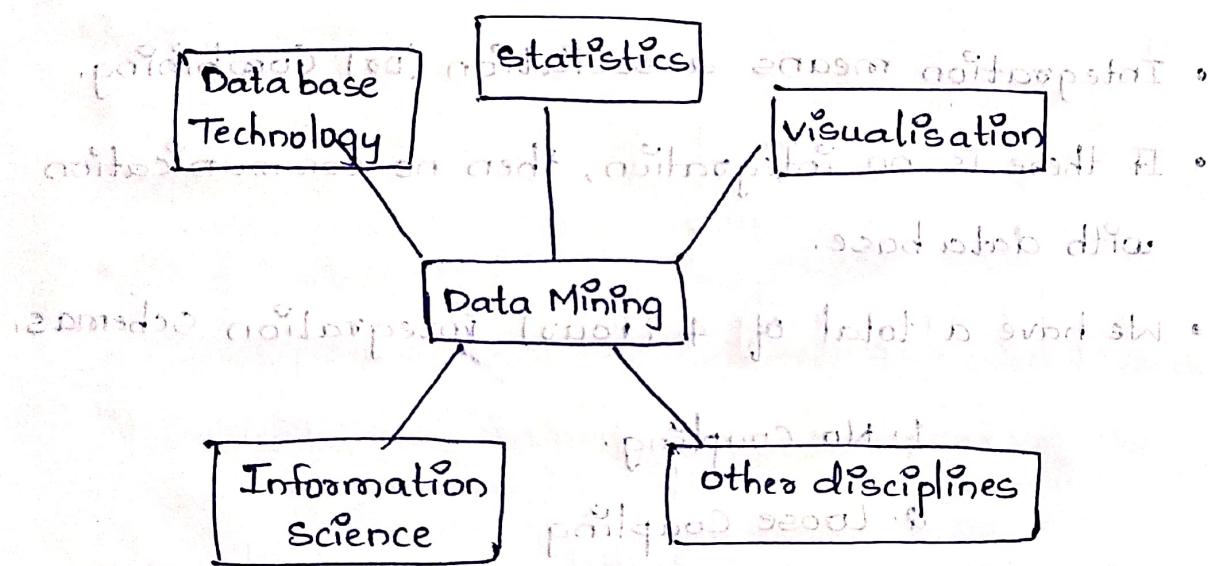
5. Classification based on techniques used:

- ML
- Statistics

- Neural Networks
- Pattern recognition
- dataware house Oriented technique, etc.

4. Classification based on application adapted:

- Finance
- Telecommunication
- Stock markets
- emails, etc.



Data Mining Task Primitives :-

- A Data Mining Task is represented in the form of a datamining query.
- Data Mining query is defined as Data mining task primitives.
- It will allow user to interactively communicate with the data mining system.
- There are 5 data mining primitives. They are :

- Set of task relevant data to be mined.
- Specifies the kind of knowledge to be mined.
- The background knowledge to be used in discovery process.
- The interestingness Measures and thresholds for pattern evaluation.
- The expected representation & visualizing the discovered patterns.

Integration of Data mining System with a Data warehouse :-

- Integration means association (or) Combining.
- If there is no integration, then no communication with database.
- We have a total of 4 (four) integration schemas.

1. No coupling.

2. Loose Coupling

3. Semitight coupling

4. Tight Coupling

1) No Coupling :-

- Data mining system will not use any function.
- That is there is no communication with Database.
- Data mining system directly communicates with other storage methods to get data.

2) Loose coupling :-

- Data Mining System will use some of the functionalities. (only upto some extent).
- It is better than No coupling.
- It is suitable for small data sets.

3) Tight Coupling :-

- Data mining system is linked to the database but not completely.
- Some of the data mining primitives are also implemented in database.

4) Tight Coupling :-

- Data mining system is completely linked to the database.
- It is most efficient among all.
- The data base system is fully integrated with such a way that it becomes a part of the data mining system.
- It is efficient and optimised implementation of data mining.

Major Issues in Data Mining :-

1. Mining different kinds of knowledge in database.
2. Interactive mining of knowledge at multiple levels of abstractions.
3. Incorporation of background knowledge.
4. Presentation and visualization of data mining results.

5. Handling noisy / incomplete data.

6. Efficiency and Scalability of data mining algorithm
(trading memory space for speed) - efficient algorithms

Data Preprocessing :-

The process of transforming raw data into an understandable format.

- 4 major tasks.

1. Data cleaning.

2. Data Integration

3. Data Reduction

4. Data Transformation.

① Data Cleaning :-

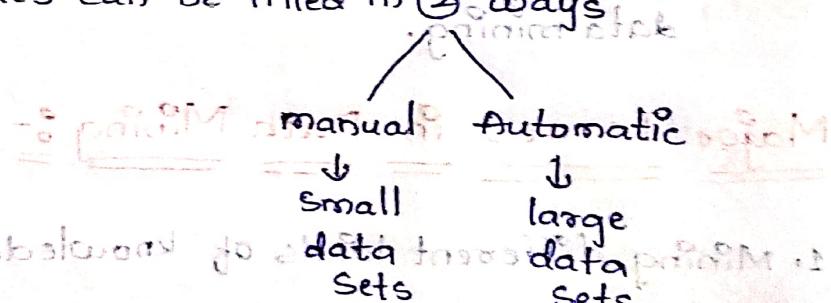
Process of removal of incorrect, incomplete, inaccurate data, also replaces missing data.

a) Handling missing values :-

- In place of missing values, we can replace with "NA", with mean values, with median values.

- Sometimes replaced with most probable values

- missing values can be filled in 2 ways



b) Handling noisy data :-

Noise data is inconsistent / error data.

Methods to handle :-

(i) Binning :-

First, data is sorted. Then sorted data is stored in bins.

③ methods to handle data in bins:

- bin - smoothing by bin mean
- smoothing by bin median
- smoothing by bin boundary.

(ii) Regression:

Numerical prediction of data.

(iii) clustering:

Similar data items are grouped at one place

dissimilar items - outside the cluster.

② Data Integration :-

Multiple heterogeneous sources of data are combined into single data set.

② types of data integration:-

a) Tight coupling :-

Data mining system is completely linked into the database.

b) Loose coupling :-

Data mining system uses some of the functionalities (only upto some extent).

③ Data Reduction :-

Volume of data is reduced to make analysis easier.

Methods :-

a) Dimensionality Reduction :-

Reduces no. of input variables in the dataset.

b) Datacube Aggregation :-

Data is combined to construct a datacube.

c) Attribute subset Selection :-

Highly relevant attributes should be used.

d) Numerosity Reduction :-

Here, we store only model of data instead of entire data.

④ Data Transformation :-

Data is transformed into appropriate form suitable for mining process.

a) Normalization :-

Done in order to scale data values in specified range (-1.0 to 1.0 or from 0 to 1).

b) Attribute selection :-

New attributes are selected by using old ones.

c) Discretization :-

Raw values are replaced by interval levels.

d) Concept Hierarchy Generation :-

Identifying attributes are converted from low level to high level.

Ex: city \rightarrow country.