

Data Mining

&

Data Warehousing

جز

Data Mining  $\Rightarrow$  The process of extracting information to identify patterns, trends, and useful data that would allow the business to take the data-driven decision from huge sets of data is called data-mining.

$\Rightarrow$  Data mining refers to the analysis of data. It is the computer supported process of analyzing huge sets of data that have either been compiled by computer systems or have been downloaded into the computer.

In the data mining process, the computer analyzes the data and extract useful information from it.

It looks for hidden patterns within the data set and try to predict future behaviour.

Data mining is primarily used to discover ~~the~~ and indicate relationships among the data sets.

$\Rightarrow$  Data Mining is one of the most useful techniques that help entrepreneurs, researchers and individuals to extract valuable information from huge sets of data.



⇒ Data mining is also called Knowledge Discovery in Database (KDD).

⇒ The Knowledge Discovery (KDD) process includes Data Cleaning, Data Integration, Data Selection, Data Transformation, Data Mining, Pattern Evaluation, and Knowledge presentation.

⇒ Data Mining is a process used by organizations to extract specific data from huge databases to solve business problems. It primarily turns raw data into useful information.

**Data Warehouse** ⇒ A data warehouse refers to a place where data can be stored for useful mining.

It is like a quick computer system with huge data storage capacity.

Data from various systems are copied to the warehouse.

⇒ Data Warehouse combines data from numerous source which ensure the data quality, accuracy and consistency.

Data flows into a data ~~base~~ warehouse from different databases.

A data warehouse work s by sorting out data into a ~~fd~~ pattern that depicts the format and types of data.

⇒ Data Warehouse and Databases both are relative data systems, but both are made to serve different purposes.

A data warehouse is used to store a huge amount of historical data and empowers fast requests over all the data, typically using Online Analytical Processing (OLAP)



⇒ A database is made to store current transactions and allow quick access to specific transactions for ongoing business process, commonly known as Online Transaction Processing (OLTP).

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

W. H. Inman, 1993

⇒ Data warehouse is non-volatile. In other words, it contains only read-only data.

⇒ Database is a collection of related data that represents some elements of the real world whereas Data Warehouse is an information system that stores historical and commutative data from single or multiple sources.

## Standard Definition of Data Warehouse

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

W. H. Inmon, 1993

**Subject-Oriented**  $\Rightarrow$  Covers all the individual applications concerning an organization.

The data warehouse is oriented around the major subjects of an enterprise, e.g. customer, vendor, product & productivity.

**Integrated**  $\Rightarrow$  The data collected in data warehouse are taken from ~~the~~ different sources. These data are integrated in a singular, globally acceptable fashion.

Data warehouse become more successful when data from multiple sources are combined in such a way that data inconsistencies and conflicts are effectively addressed. This process is known as data scrubbing.

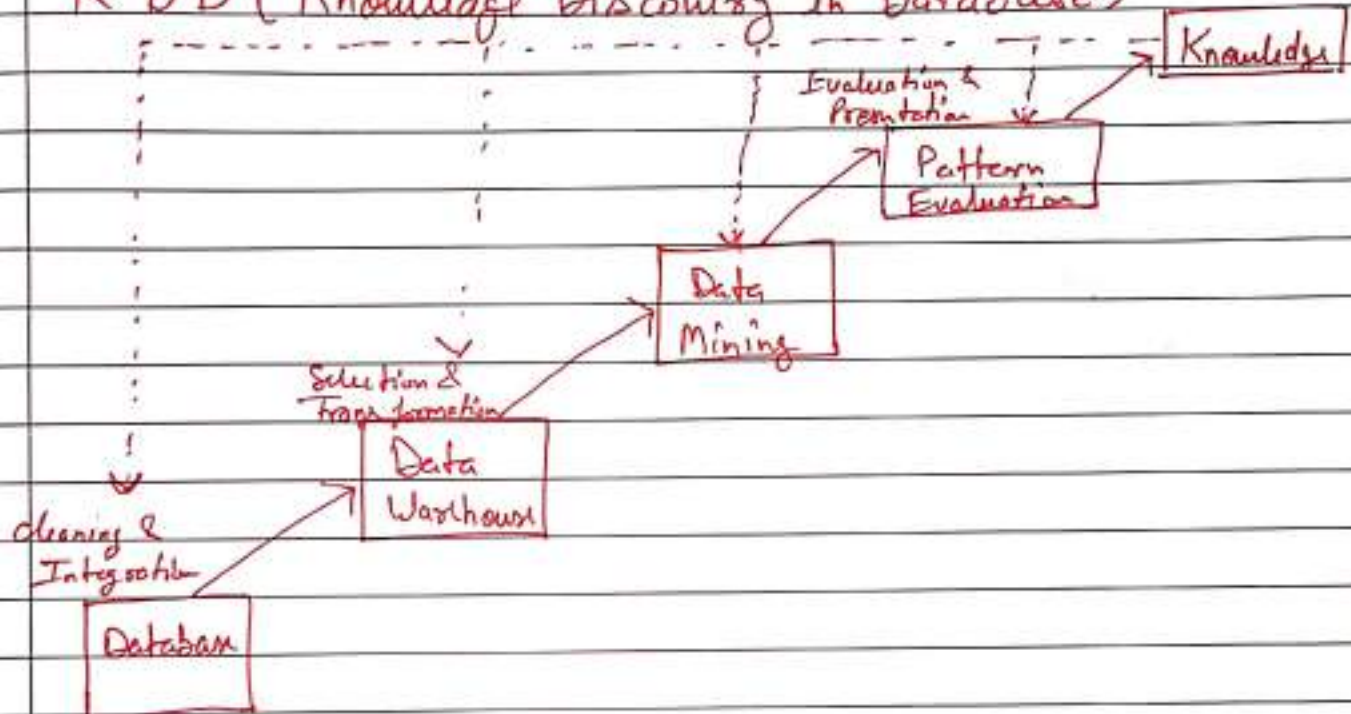
**Time-Variant**  $\Rightarrow$  The data in a data-warehouse is time-dependent. The data warehouse focus on changes over time.



Non-Volatile  $\Rightarrow$  The data warehouse is a non-volatile. In other words, it contains read-only data.

Once the data is loaded into the warehouse from an application-oriented environment or external sources, it does not change and is only available for accessing.

## KDD (Knowledge Discovery in Database)



**Data Cleaning** ⇒ It is the process of removing noisy or inconsistent data.

⇒ Cleaning is performed for detection of syntax errors.

⇒ Parser decides whether the given string of data is acceptable within data specification.

**Data Integration** ⇒ Different Data from multiple sources are combined.

**Data Selection** ⇒ Useful and relevant data are retrieved from the database for analysis.



Data Transformation  $\Rightarrow$  Data are transformed or consolidated into forms appropriate for mining by performing summary or aggregation operations, for instance.

Data Mining  $\Rightarrow$  It is a process where intelligent methods are applied in order to extract data patterns.

Pattern Evaluation  $\Rightarrow$  To identify the truly interesting patterns representing knowledge-base on some specific criteria.

Knowledge Representation  $\Rightarrow$  Where visualization and knowledge representation techniques are used to present the mined knowledge.

## Types of Data Mining

Data mining can be performed on the following types of data -

- ① Relational Database
- ② Data Warehouse
- ③ Data Repositories
- ④ Object-Relational Database
- ⑤ Transactional Database

## Advantages of Data Mining

## Disadvantages of Data Mining

## Applications of Data Mining

These are the following areas where data mining is widely used -

- ① Data Mining in Healthcare
- ② Data Mining in Market Based Analysis
- ③ Data Mining in education
- ④ Data Mining in Manufacturing engineering
- ⑤ Data Mining in Customer Relationship Management (CRM)
- ⑥ Data Mining in Fraud Detection
- ⑦ Data Mining in Lic Detection.
- ⑧ Data Mining in Financial Banking.



## Challenges in Implementation of Data Mining.

- ① Incomplete and noisy data.
- ② Data Distribution
- ③ Complex data.
- ④ Performance
- ⑤ Data Privacy & Security
- ⑥ Data Visualization

## Data Mining Techniques

There are some important data mining techniques -

- ① Classification
- ② Clustering
- ③ Regression
- ④ Association Rules
- ⑤ Outlier Detection
- ⑥ Sequential Patterns
- ⑦ Prediction