

UNIT-1 Introduction to Data Mining

Related Technologies-Machine Learning, DBMS, OLAP and Statistics, Data Mining Goals, Stages of Data Mining Process, Data Warehouse and Multidimensional Data Model, OLAP Operations.

UNIT-2 Data Preprocessing and Data Visualization

Data Cleaning, Data Transformation, Data Reduction, Discretization and Concept hierarchy, Exploratory Data Analysis Tools(Plots, Graphs, Summary Statistics, histograms, heat maps)

UNIT-3 Data Mining algorithms:-Association Rules and Regression Analysis

Item sets, Frequent Patterns, Interestingness measures(support, confidence and lift), Correlation analysis,

Apriori and Frequent Growth Algorithms, Linear Regression and Models

10

UNIT-4 Data Mining Algorithms: Classification

Decision Trees, Random Forests, Bayesian Networks, Nearest Neighbour Algorithms.

UNIT-5 Data Mining Algorithms:Clustering and Model Evaluation:

Partitioning Methods:-K-means,k-Medoids,Expectation Maximization,Hierarchical Methods:Distance-Based agglomerative clustering and divisive clustering,Training And Testing(Cross-Validsation),Combining multiple models(bagging,boosting)

① Data mining is the extraction of interesting patterns or knowledge from a huge amount of data

Data mining can be applied to any kind of data as long as the data are meaningful for target application

Stages of Data mining

① **Data Gathering** : Relevant data for analytical application is assembled
(Data fetching)

from various sources {mix of Structured & unstructured data}

② **Data preparation** ; This stage includes set of steps to ready the data for mining which includes;

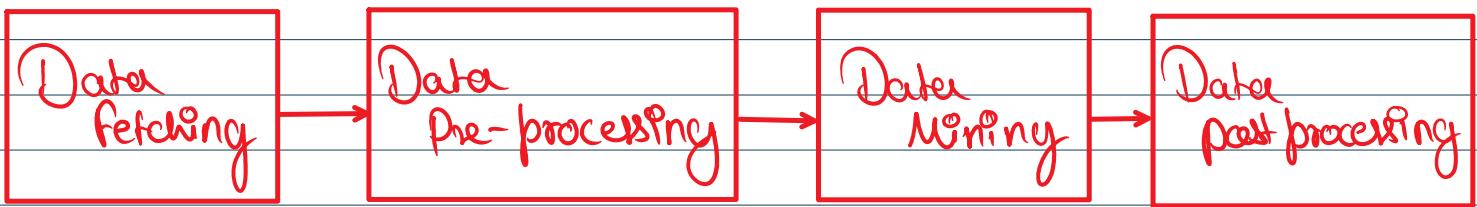
- data Exploration
- profiling
- pre processing
- data cleaning
- data transformation

③ **Mining of data** ; Prepared data is mined according to appropriate mining techniques chosen by data scientist
(Data mining)

Various techniques / algorithms are applied

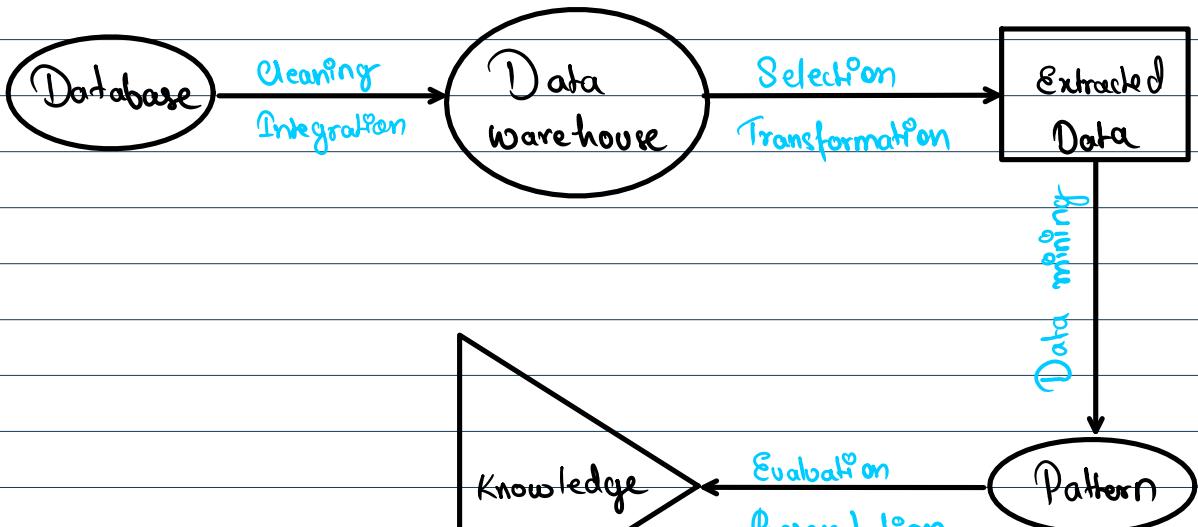
④ **Data Analysis & Interpretation**
(Data post processing)

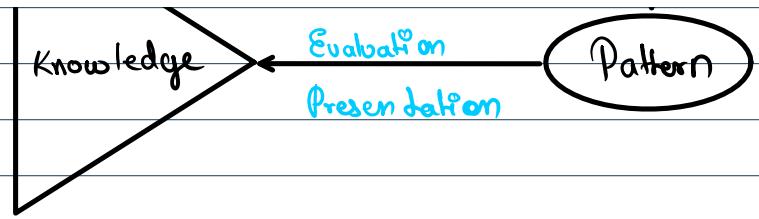
Mining results are used to create analytical models
(which helps in business decision making)



Data mining includes following iterative steps

- ① Data Cleaning (removal of noise & inconsistent data)
- ② Data Integration (Combination or Combination of data from multiple data sources)
- ③ Data Selection (selection of relevant data for analysis from DB)
- ④ Data transformation (Data is converted into forms appropriate for mining)
- ⑤ Data mining (essential process where intelligent methods are applied to extract data patterns)
- ⑥ Pattern Evaluation (Identification of pattern representing knowledge)
- ⑦ Knowledge presentation (Visualization of knowledge & their techniques.)





Multi-Dimensional View of Data Mining

■ Data to be mined ↗ *Data*

- Database data (extended-relational, object-oriented, heterogeneous, legacy), data warehouse, transactional data, stream, spatiotemporal, time-series, sequence, text and web, multi-media, graphs & social and information networks

■ Knowledge to be mined (or: Data mining functions) ↗ *Knowledge*

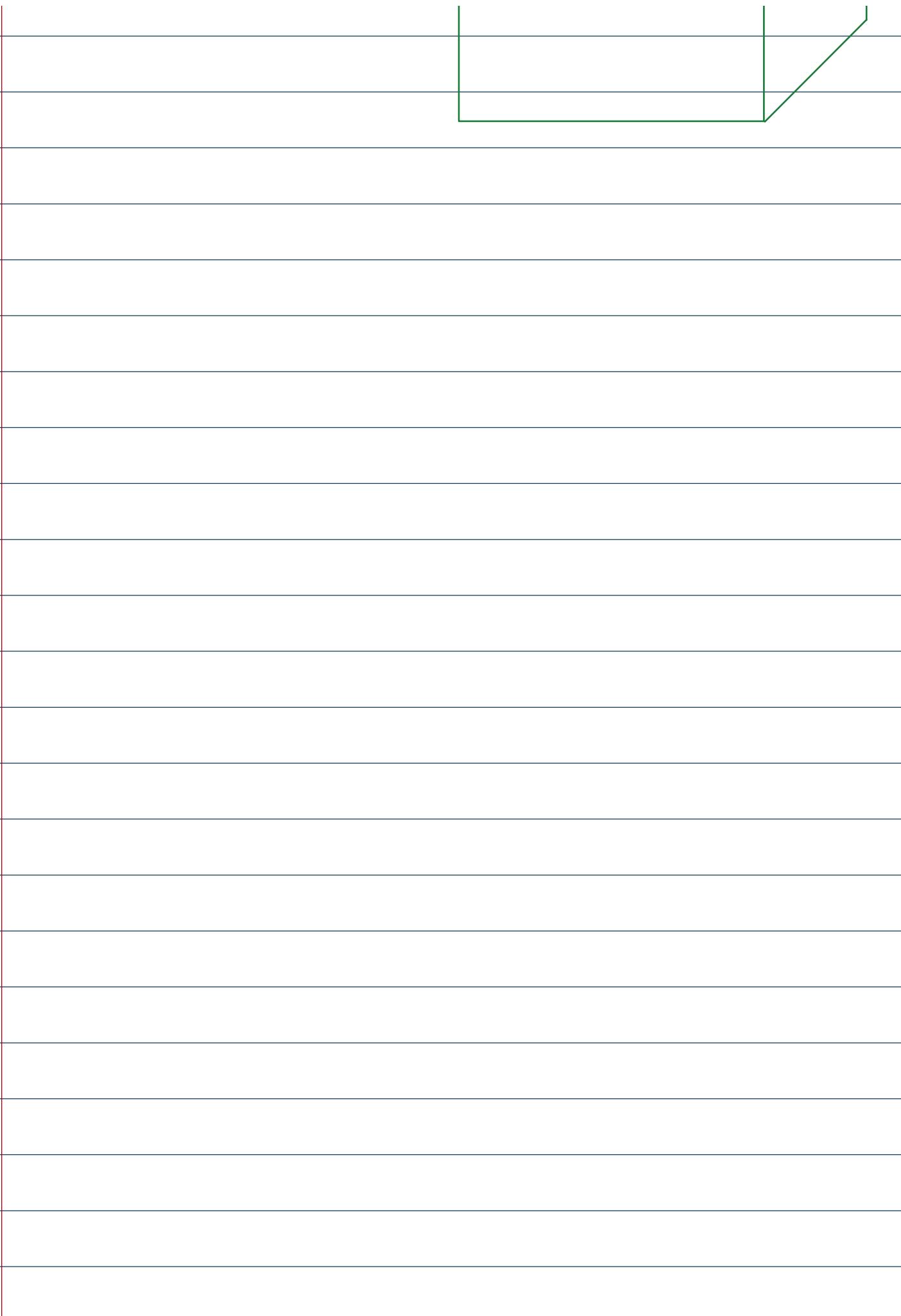
- Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
- Descriptive vs. predictive data mining
- Multiple/integrated functions and mining at multiple levels

■ Techniques utilized ↗ *Technology*

- Data-intensive, data warehouse (OLAP), machine learning, statistics, pattern recognition, visualization, high-performance, etc.

■ Applications adapted ↗ *Application*

- Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.

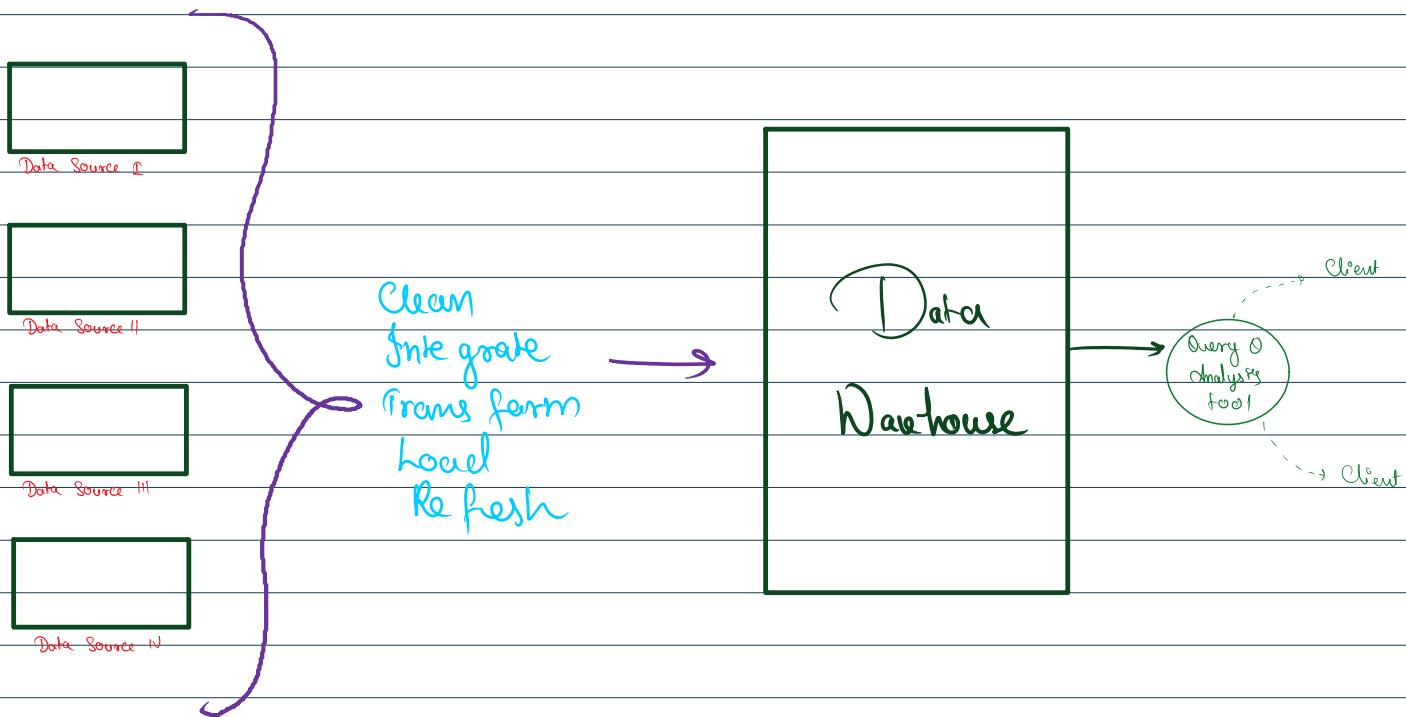


Data Warehouse ; It is a repository for long-term storage of data from multiple sources

- The data is stored under a unified schema
 - & typically summarized
- Data warehousing allows multidimensional data analysis capability



Online analytical process (OLAP)



OLAP

- Consists of long term historical data collected from various databases
- Makes use of data warehouse
- The data is used in long term planning & decision making
- Astronomical size of data
- Not often updated
- Managed by CEO, GM
- Only read operation is required
- Snowflake or Star Schema

OLTP

- Consists of short term, currently operational data & transactional data
- Makes use of Standard DBMS
- Data is used to perform day-to-day fundamental operations
- Relatively small size of data
- data integrity & files are constantly being updated
- Managed by clerks, Managers
- Both read & write required
- ER diagram

Types of Attributes

① Nominal attr.

② Binary

③ Ordinal

④ Nomeric

!Types of Data Sets

28 September 2022 07:14 AM

1 Data Sets

Schema Diagram

Star Schema

Time key
day
week
Month
years

Item key
Name
brand
type

fact table

Branch key
Name
type

Time key
Item key
branch key
Item key
Units Sold
\$ Sold

Item key
Sheet
Copy
State
Country

Other than the fact table rest are called the dimension table

If any dimension table were to be normalized then

Advantages of Decision trees

- ① It doesn't require any domain knowledge
- ② Classification steps of decision tree are simple & fast
- ③ Missing values in data doesn't affect output
- ④ A decision tree model is automatic & doesn't need Standardization of data