# UNIT 1 - FUNDAMENTAL CONCEPTS OF DATA MINING

*"The world is **data** rich but **information** poor. Data mining tools that can turn data tombs into "golden nuggets" of knowledge."*

**(1) Define : Data Mining**

- Data mining is the process of discovering interesting patterns and knowledge from large amounts of data.

**(2) Data Warehouse :**

- Data repository architecture
- This is a repository of multiple heterogeneous data sources organized under a unified schema at a single site to facilitate management decision making.
- Data warehouse technology includes data cleaning, data integration, and online analytical processing (OLAP)
  - o OLAP is analysis techniques with functionalities such as summarization, consolidation, and aggregation
  - o Although OLAP tools support multidimensional analysis and decision making, additional data analysis tools are required for in-depth analysis - for example, data mining tools that provide data classification, clustering, outlier/anomaly detection

**(3) Difference between Database and Data Mining:**

|  | Database | Data Mining |
|---|---|---|
| **Purpose** | A database is a structured system for storing, managing, and retrieving raw data efficiently | Data Mining is the analytical process of discovering hidden patterns, trends, and insights from large datasets to make predictions and decisions |

| | Data management, efficiency, reliability, and data integrity | Data analysis, pattern recognition, prediction, and knowledge discovery |
|---|---|---|
| **Focus** | | |
| **Data Type** | Transactional, current, detailed data | Large volumes of historical, aggregated data |
| **Example** | SQL Database, NoSQL Database | Market Basket Analysis, Fraud Detection |

## (4) Common terms used: patterns, trends, & predictions

- ➢ Patterns:
    - ▪ A pattern is a sequence of data points that repeats in a recognizable, predictable way with a consistent structure.
    - ▪ Unlike a general trend, a pattern is cyclical or recurring.
    - ▪ Example: Increased retail sales during the holiday season or higher ice cream sales in the summer
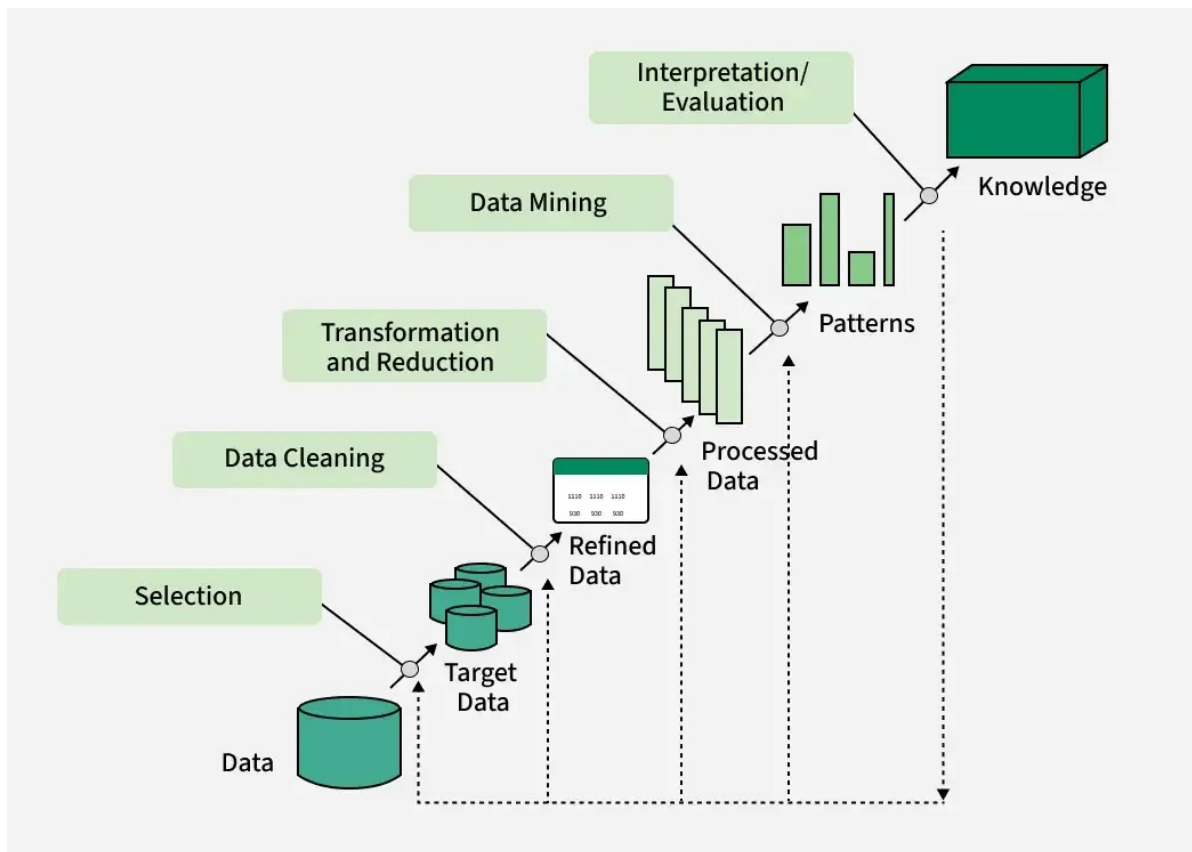- ➢ Trends:
    - ▪ A trend is the general, long-term direction in which a variable or market is moving over a period of time.
    - ▪ Trends describe the overall movement, even if there are short-term fluctuations or patterns along the way
    - ▪ Example: Uptrend (Bull Market) or Downtrend (Bear Market)
- ➢ Predictions:
    - ▪ The systematic process of using data analysis techniques to project probable future events, outcomes or scenarios
    - ▪ Example: Marketing manager predicting how much a specific customer will spend during a sale, using past customer data like purchase history, demographics etc.

## (5) Steps in Knowledge Discovery Process (KDD)

- Knowledge Discovery in Databases (KDD) refers to the complete process of uncovering valuable knowledge from large datasets.



1. Data Selection:
   - Focusing on what's needed for the specific analysis
   - Identifying and retrieving relevant data from various sources
2. Data Preprocessing:
   2.1) Data Cleaning:
      - Removing noise
      - Handling missing values
      - Removing duplicates
   2.2) Data Integration:
      - Combining data from multiple sources to create a consistent dataset
3. Data Transformation:

- Converting data into formats suitable for mining, often involving aggregation or normalization

4. Data Mining:

- Applying intelligent algorithms (like clustering, classification, association rules) to discover patterns and relationships

5. Pattern Evaluation:

- Assessing the discovered patterns to determine their interestingness, novelty, and usefulness

6. Knowledge Representation:

- Presenting the final, validated knowledge in an understandable form for users, often using visualization techniques like reports or charts

## (6) Importance of Data Mining

- Data mining ensures that useful information can be derived from raw data and used to benefit both the organization and its customers.
- Data mining ensures data-driven decision-making
- Data mining helps in analyzing substantial amounts of data quickly
- Data mining helps in identifying patterns and trends and detecting fraud
- Businesses can get reliable information through data mining
- Data mining is a cost-effective and efficient option.

## (7) Applications of Data Mining in real life

### 7.1) Education: analyzing student performance

- Predictive Analytics: Identifying students likely to struggle or drop out early, allowing timely intervention.
- Personalized Learning: Creating customized learning experiences and recommendations based on individual needs, styles, and performance trends.
- Early Intervention: Detecting patterns indicating disengagement or difficulty in specific subjects for focused support.

- Curriculum Optimization: Analyzing which teaching methods and course materials yield the best results to refine curriculum design.
- Resource Allocation: Informing decisions on where to best allocate support staff, tutoring, or funding for maximum impact.
- Performance Prediction: Forecasting final exam results based on internal assessments and student behavior.

## 7.2) Business: customer purchase analysis

- Market Basket Analysis: Identifies items frequently bought together (e.g., bread & butter) to inform product placement, cross-selling, and strategic bundling deals.
- Personalized Recommendations: Analyzes browsing/purchase history to suggest relevant products, increasing conversion rates and customer satisfaction (e.g., Amazon, Netflix).
- Customer Segmentation: Groups customers by behavior (e.g., frequent buyers, bargain hunters, seasonal shoppers) to create specific, effective marketing campaigns.
- Customer Retention & Churn Prediction: Predicts which customers are likely to leave (churn) based on patterns, allowing businesses to proactively offer incentives to retain them.
- Demand Forecasting: Uses historical data and trends to predict future product demand, optimizing inventory and reducing stockouts or overstock.
- Targeted Marketing: Develops highly customized promotions and loyalty programs for specific customer segments
- Fraud Detection: Identifies unusual purchase patterns that might indicate fraudulent activity, protecting both the business and customers.

## 7.3) Healthcare: disease prediction

- Early Diagnosis: Identifying subtle patterns in medical images (X-rays, MRIs) or lab results (blood tests, vitals) for early detection of tumors, stroke, or chronic diseases like diabetes.

- Risk Assessment: Predicting a patient's likelihood of developing conditions like heart disease by analyzing factors such as age, blood pressure, and lifestyle data.

- Prognosis & Survival: Building models (e.g., using Random Forest) to forecast disease outcomes and survival rates for patients with advanced illnesses like cancer.

- Clinical Decision Support: Helping doctors make faster, more informed diagnoses by rapidly processing patient data and highlighting potential diagnoses or complications.

- Personalized Treatment: Classifying patients into groups to find the most effective treatment protocols for specific profiles, leading to customized care plans.

## 7.4) Banking: fraud detection

- Anomaly Detection: Identifies transactions that deviate from a customer's normal behavior, such as sudden large purchases in distant locations or unusual transaction times.

- Pattern Recognition: Uncovers common characteristics of past fraud to build models that predict future fraudulent attempts.

- Clustering: Groups similar transactions or accounts, allowing analysts to identify clusters of potentially fraudulent activity.

- Classification: Uses algorithms (like Neural Networks, SVM) to categorize transactions as legitimate or fraudulent.