# Data Mining and Its Benefits

**Data mining**, also known as knowledge discovery in data (KDD), is the process of uncovering patterns and valuable information from large data sets.

- Evolution of data warehousing technology and the growth of big data have accelerated the adoption of data mining techniques.
- Companies transform raw data into useful knowledge.
- Despite evolving technology, challenges with scalability and automation persist.

Data mining has improved organizational decision-making through insightful data analyses. The techniques can either describe the target dataset or predict outcomes using machine learning algorithms.

**Purposes of data mining:**

- Predicting various outcomes
- Modeling target audience
- Collecting information about products

Data mining analyzes data and converts it into meaningful information, helping businesses make accurate and better decisions. It aids in:

- Developing smart market decisions
- Running accurate campaigns
- Making predictions
- Analyzing customer behaviors and insights

**Benefits of Data Mining:**

- Helps companies gather reliable information
- Efficient, cost-effective solution compared to other data applications
- Aids in making profitable production and operational adjustments
- Utilizes both new and legacy systems
- Facilitates informed decision-making
- Detects credit risks and fraud
- Allows data scientists to analyze large amounts of data quickly
- Detects fraud, builds risk models, and improves product safety

- Initiates automated predictions of behaviors and trends
- Discovers hidden patterns

# Types of Data that can be Mined

Data mining is applicable to various types of media and data repositories. Algorithms and approaches may differ based on data type.

**Applicable data types:**

- Relational databases
- Object-relational and object-oriented databases
- Data warehouses
- Transactional databases
- Unstructured and semi-structured repositories (e.g., the Web, social media)
- Advanced databases (e.g., spatial, multimedia, time-series, textual databases)
- Flat files

# How Data Mining Works

**Cross-Industry Standard Process for Data Mining (CRISP-DM)** is a methodology guiding data mining efforts. It includes descriptions of typical project phases, tasks involved, and task relationships. The CRISP-DM model is flexible and customizable.

**Phases of CRISP-DM:**

1. **Business Understanding**
   - Define objectives and problems.
   - Determine required data.
2. **Data Understanding**
   - Collect relevant data from various sources.
   - Ensure data encompasses necessary datasets.
3. **Data Preparation**
   - Extract, transform, and load (ETL) data.
   - Clean, populate null sets, remove duplicates, resolve errors, and allocate data into tables.
4. **Modeling**
   - Address relevant data sets and select statistical/mathematical approaches.
   - Use techniques like classification, clustering, and regression analysis.
5. **Evaluation**
   - Evaluate model efficiency in answering business questions.
   - Adjust models or data as needed.
6. **Deployment**
   - Implement models via visual presentations, reports, or actionable strategies (e.g., new sales strategies, risk-reduction measures).

# Classification of Data Mining Systems

Data mining systems can be categorized based on various criteria:

- **Data Source Mined**: Spatial data, multimedia data, time-series data, text data, Web data, etc.
- **Data Model Drawn On**: Relational databases, object-oriented databases, data warehouses, transactional data, etc.
- **Knowledge Discovered**: Characterization, discrimination, association, classification, clustering, etc.
- **Mining Techniques Used**: Machine learning, neural networks, genetic algorithms, statistics, visualization, etc.
- **User Interaction**: Query-driven, interactive exploratory, or autonomous systems.

# Data Mining Techniques

Data mining is used to identify patterns and derive business insights. Common techniques include:

1. **Classification Analysis**
   - Assign data points to groups based on specific questions/problems.
2. **Association Rule Learning**
   - Uncover relationships between data points (e.g., business travelers' room choices and dining habits).
3. **Anomaly or Outlier Detection**
   - Find unusual data within a set to detect fraud or unusual sales patterns.
4. **Clustering Analysis**
   - Separate data points with common traits into subsets (e.g., customer segmentation).
5. **Regression Analysis**
   - Understand important factors and their interactions within a dataset (e.g., predicting sales based on weather forecasts).

# Data Mining vs Data Warehousing

| Attribute | Data Mining | Data Warehouse |
|---|---|---|
| **Process** | Analyzing unknown patterns of data. | Database system designed for analytical instead of transactional work. |
| **Method** | Comparing large amounts of data to find patterns. | Centralizing data from different sources into one common repository. |
| **Users** | Usually done by business users with assistance of engineers. | Occurs before any data mining can take place. |
| **Purpose** | Extracting data from large data sets. | Pooling all relevant data together. |
| **Benefits** | Detection and identification of errors in the system. | Consistent updates; ideal for business owners seeking the latest features. |

| Attribute | Data Mining | Data Warehouse |
| --- | --- | --- |
| Applications | Creating suggestive patterns of important factors like buying habits, products, sales. | Adds value to operational business systems like CRM systems. |
| Accuracy | Techniques are never 100% accurate and may have serious consequences. | Data required for analysis may not be integrated, leading to loss of information. |
| Misuse | Information gathered can be misused against a group of people. | High maintenance system impacting revenue of medium to small-scale organizations. |
| Workload | Complicated queries increase workload. | Complicated to implement and maintain. |
| Analytical Tool | Equips organizations with pertinent and usable knowledge-based information. | Stores a large amount of historical data for analyzing different time periods and trends. |
| Resources | Requires significant resources for training and implementation. | Data pooled from multiple sources needs to be cleaned and transformed, which is challenging. |
| Cost-Effectiveness | Cost-effective and efficient compared to other statistical data applications. | Simplifies business data; users primarily input raw data. |
| Error Identification | Identifies errors that can lead to losses. | Allows access to critical data from multiple sources in one place, saving retrieval time. |
| Actionable Strategies | Helps generate actionable strategies built on data insights. | Once data is input, it's easy to retrieve and unlikely to be lost. |

# Data Mining Tools

Data mining techniques utilize domain knowledge from statistical analysis, artificial intelligence, and database systems to analyze data in different dimensions and perspectives. They discover patterns or trends from large data sets and transform data into useful information for decision-making.

## Popular Data Mining Tools

- **Orange**
  - Developed at the bioinformatics lab, Ljubljana University, Slovenia.
  - Machine learning and data mining software suite.
  - Supports data visualization and is written in Python.
  - Uses components called "widgets" for various functionalities:
    - Displaying data table and allowing feature selection.
    - Reading data.
    - Training predictors and comparing learning algorithms.
    - Visualizing data elements.
  - Can perform data mining via visual programming or Python scripting.

- Supports various visual tools: bar charts, scatter-plots, trees, dendrograms, heat maps.
- Machine learning components, bioinformatics, and text mining add-ons.
- **SAS Data Mining**
  - Stands for Statistical Analysis System.
  - Product of SAS Institute for analytics and data management.
  - Mines data, manages information, and analyzes statistics.
  - Graphical UI for non-technical users.
  - Distributed memory processing architecture, highly scalable.
- **Rattle Data Mining**
  - GUI for data mining using R.
  - Presents statistical and visual summaries of data.
  - Builds both unsupervised and supervised machine learning models.
  - Free open-source software.
  - Captures interactions as an R script.
- **Rapid Miner**
  - Popular tool for predictions, written in JAVA.
  - Integrated environment for text mining, deep learning, machine learning, predictive analysis.
  - Supports applications in commercial, research, education, and more.
  - Client/server model, template-based frameworks for fast delivery.
- **Data Melt**
  - Free tool for numeric computation, mathematics, data analysis, data visualization.
  - Supports scripting languages like Python, Ruby, Groovy.
  - Offers statistics, analysis of large data volumes, scientific visualization.
  - Creates high-quality vector-graphics images.
  - Faster than standard Python implemented in C.

# Applications of Data Mining

Data mining is primarily used by companies in retail, finance, communication, and marketing to analyze transactional data and determine pricing, customer preferences, product positioning, impact on sales, customer satisfaction, and profits.

- **Basket Analysis**
  - Analyzes consumer buying habits to recommend purchases.
  - Used in retail and law enforcement.
- **Sales Forecasting**
  - Predictive analysis to project sales and set targets.
  - Examines historical data, financial indicators, consumer spending habits, and trends.
- **Bioinformatics**
  - Extracts useful knowledge from biological data.
  - Applications include gene finding, disease diagnosis, treatment optimization, data cleansing.
- **Inventory Planning**
  - Provides up-to-date information on product inventory, delivery schedules, production requirements.
  - Helps manage product stock efficiently.
- **Customer Segmentation**

- Segments customers to increase market effectiveness.
    - Helps retain customers with tailored offers.
- **Customer Relationship Management (CRM)**
    - Improves customer loyalty and implements customer-focused strategies.
    - Uses data mining for analysis of collected data.
- **Healthcare**
    - Predicts patient volume, ensures appropriate care, detects fraud, and abuse.
- **Education**
    - Develops methods to discover knowledge from educational data.
    - Predicts student behavior and results, improving teaching methods.
- **Intrusion Detection**
    - Improves detection by focusing on anomaly detection.
    - Extracts relevant data for analysis.
- **Criminal Investigation**
    - Identifies crime characteristics and relationships.
    - Converts crime reports into word processing files for analysis.
- **Fraud Detection**
    - Observes data anomalies to detect fraud.
    - Adopted by banking, financial institutions, and SaaS-based companies.
- **Operational Optimization**
    - Reduces costs across operational functions.
    - Identifies bottlenecks and improves decision-making.

# Issues in Data Mining

Data mining involves complex algorithms and data integration from various sources. Major issues include:

## Mining Methodology and User Interaction

- **Mining different kinds of knowledge**
    - Different users require different kinds of knowledge.
- **Interactive mining**
    - Allows users to focus and refine data mining requests.
- **Incorporation of background knowledge**
    - Guides discovery and expresses discovered patterns.
- **Data mining query languages**
    - Should be integrated with data warehouse query language for efficient mining.
- **Presentation and visualization**
    - Patterns need to be expressed in understandable high-level languages and visual representations.
- **Handling noisy/incomplete data**
    - Data cleaning methods are necessary for accurate pattern discovery.
- **Pattern evaluation**
    - Patterns should be interesting, representing common knowledge or novelty.

# Performance Issues

- **Efficiency and scalability**
  - Algorithms must be efficient and scalable for large databases.
- **Parallel, distributed, and incremental mining**
  - Algorithms should handle huge databases and distribute data for parallel processing.

# Diverse Data Types Issues

- **Handling complex data types**
  - Systems must manage complex data objects like multimedia, spatial, and temporal data.
- **Mining from heterogeneous databases**
  - Challenges arise from structured, semi-structured, and unstructured data sources.

# Check Your Progress-1

1. What are the descriptive and predictive data mining techniques? Explain.

# Check Your Progress-2

1. Identify and describe the features of some more open source / free data mining tools which were not discussed in this course content.