**Master of Computer Application**

# MCA-E6N
# Data Mining

उ० प्र० राजर्षि टण्डन
मुक्त विश्वविद्यालय, प्रयागराज

# Block
# 1

## Data Pre-processing and Data Warehousing

## Course Design Committee

| | |
|---|---|
| **Prof. Ashutosh Gupta** | Chairman |
| Director (In-charge) | |
| School of Computer & Information Science, UPRTOU Allahabad | |
| **Prof. Suneeta Agarwal** | Member |
| Dept. of Computer Science & Engineering | |
| Motilal Nehru National Institute of Technology Allahabad | |
| **Dr. Upendra Nath Tripathi** | Member |
| Associate Professor | |
| DeenDayalUpadhyay Gorakhpur University, Gorakhpur | |
| **Dr. Ashish Khare** | Member |
| Associate Professor | |
| Dept. of Computer Science, University of Allahabad, Prayagraj | |
| **Ms. Marisha** | Member |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |
| **Mr. Manoj Kumar Balwant** | Member |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |

## Course Preparation Committee

| | |
|---|---|
| **Dr. Tulika Narang** | (Block 1, 2 & 4) Author |
| Assistant Professor, Computer science | |
| United University, Rawatpur, Prayagraj | |
| **Dr. Krishan Kumar** | (Block 3-Unit 7) Author |
| Assistant Professor, | |
| Department of Computer Science, Faculty of Technology | |
| Gurukula Kangri Vishwavidyalaya, Haridwar (UK) | |
| **Dr. Pooja Yadav** | (Block 3-Unit 8) Author |
| Assistant, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Iram Naim** | (Block 3-Unit 9) Author |
| Assistant Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Brajesh Kumar** | Editor |
| Associate Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Mr. Manoj Kumar Balwant** | Coordinator |
| Assistant Professor (Computer Science), | |
| School of Sciences, UPRTOU, Prayagraj | |

# Block-1 Introduction

This block provides a comprehensive exploration of data types, data pre-processing, and data warehousing concepts. The first unit focuses on the fundamental data types relevant to the data mining process. It explores various data pre-processing techniques essential for preparing data for analysis. Additionally, data visualization is highlighted as a critical component for interpreting data effectively. The second unit explores the implementation of data warehouses, discussing their necessity and core characteristics. This unit also clarifies the distinction between data warehouses and data marts, providing insight into their respective roles.

Furthermore, we examine the client-server computing model and its application in data warehouse environments. The third unit introduces critical components of data warehouse management, including the load manager, warehouse manager, query manager, and end-user tools. Here, we will discuss the three-layered architecture of a data warehouse, which encompasses the database, analytical processing modules, and front-end tools. The unit addresses various implementation challenges, detailing the hardware and software requirements for a successful data warehouse deployment. Comprehensive definitions of a data warehouse, its components, and its layered architecture are provided to enhance understanding.

# UNIT 1: **Processing and Visualization of Data**

**Structure**

## 1.0 Introduction

The unit focuses on data mining and related essential concepts. Data mining is the process of extracting hidden and valuable patterns of knowledge from data. Data mining includes access to data and data preparation for mining, analysis and interpretation of results. Data mining refers to the process of discovering patterns, relationships, and insights from large sets of data. It involves using various techniques and algorithms to extract valuable information from structured, semi-structured, and unstructured data sources. Data Mining builds models to identify valuable patterns in data. The input to the Data Mining process is data of various types that is prepared for the mining process. Finally the results are analyzed with data visualization techniques such as tables and graphs. In this unit, the emphasis is to understand the essential data types in data mining process. Various data pre-processing methods are also discussed in this unit. Data visualization is also an essential component of this unit.

## 1.1 Objectives

After the end of this unit, you should be able to:

- Explain different types of data used in Data Mining process.

- Understand data-preprocessing and various methods of data-preprocessing.

- Understand and able to implement different measures of similarity.

- Understand data visualization and its advantages.

# 1.2 Data Types

In data mining, various data types are used to represent and analyze the data. The choice of data type depends on the nature of the data and the specific requirements of the mining task. Here are some common data types in data mining:

❖ Numeric Data: This type of data consists of numerical values and is often used in quantitative analysis. Numeric data can be further categorized as continuous or discrete. Continuous data includes measurements such as height, weight, temperature, etc., which can take any real value within a certain range. Discrete data includes whole numbers or counts, such as the number of items sold, number of customers, etc.

❖ Categorical Data: Categorical data represents variables that can take on a limited number of distinct values or categories. Examples include gender (male/female), color (red/blue/green), occupation (doctor/engineer/teacher), etc. Categorical data can be further divided into nominal and ordinal data types. Nominal data has categories without any inherent order or ranking, while ordinal data has categories with a specific order or ranking.

❖ Text Data: Text data represents unstructured textual information, such as documents, emails, social media posts, etc. Text mining techniques are used to extract meaningful insights from this type of data. Preprocessing steps like tokenization, stemming, and entity extraction are typically performed to transform text data into a more structured format suitable for analysis.

❖ Time-Series Data: Time-series data represents measurements taken at different points in time and is commonly used for analyzing trends, patterns, and forecasting. It is characterized by the sequential nature of data points and the dependence on the temporal ordering. Examples include stock prices, weather data, sensor readings, etc.

❖ Spatial Data: Spatial data represents information related to geographic or spatial locations. It includes coordinates, shapes, distances, and other spatial attributes. Spatial data mining techniques are used to analyze patterns and relationships in spatial data. Examples include maps, satellite images, GPS coordinates, etc.

❖ Multimedia Data: Multimedia data encompasses various types of non-textual data such as images, videos, audio files, and other multimedia content. Data mining techniques are employed to extract patterns, features, and information from multimedia data. This field is also known as multimedia mining or multimedia content analysis.

❖ Graph Data: Graph data represents entities and their relationships using nodes and edges. Graph mining techniques are used to analyze the structure and properties of complex networks. Examples include social networks, web graphs, biological networks, etc.

The essential data types in data mining process can also be categorised into four main types:

❖ Nominal data.

❖ Ordinal data.

❖ Discrete data.

❖ Continuous data

Nominal data is a type of qualitative information. It labels the variables without providing the numerical value. It cannot be ordered and measured. Some examples of nominal data are letters and gender. Ordinal data follows a natural order. The essential characteristic of the nominal data is that the difference between the data values is undetermined. This type of data is found in surveys, finance and surveys. Discrete data considers discrete values. It is type of quantitative data. Examples include number of students in class. Continuous data considers data that has infinite number of probable values within a specific range. Example include Temperature.

---

*Check your progress 1*

*1. Discuss essential types of data used in data mining process.*

*2. Differentiate between ordinal, nominal, discrete and continuous data.*

---

## 1.3 Data Preprocessing

Data-preprocessing is an essential element of data preparation. It is related to converting raw data into a format useful for data preprocessing. It is an important preliminary phase for the data mining process. This involves collecting, cleaning, integrating, and transforming the raw data into a format suitable for analysis. It may include tasks like removing duplicates, handling missing values, and resolving inconsistencies.

It is concerned with the data cleaning, data transformation and data integration to prepare data for processing and analysis. The objective of data pre-processing is to improve the quality of the data. Also, it makes data more appropriate for particular data mining task.

# 1. Data Cleaning:

The process of data cleaning deals with handling irrelevant and missing values. There are various methods for performing data cleaning operation as discussed here.

**(a) Dealing with missing values:**

The following ways are used to deal with missing values-

1. Ignore the tuples that contain missing values  and removing the missing values. This can be done in following ways:

- **Listwise Deletion**: Remove any tuples with missing values. This is straightforward but can lead to a significant loss of data, especially if many tuples have at least one missing value.
- **Pairwise Deletion**: Only exclude the missing data points from the analysis, allowing the rest of the tuple to be used. This can be useful if missing data is sparse and randomly distributed.

2. Imputation to replace missing values with other values in the following ways:

- **Mean/Median/Mode Imputation**: Replace missing values with the mean (for numerical data), median, or mode (for categorical data) of the attribute. This is simple but can reduce variance and distort relationships.
- **K-Nearest Neighbors (KNN) Imputation**: Use the average of the k-nearest neighbors' values to impute the missing data. This method can be more accurate than mean imputation but is computationally expensive.
- **Regression Imputation**: Predict the missing value using a regression model built from the non-missing values.
- **Multiple Imputation**: Generate multiple imputations to account for the uncertainty in the missing values, often by creating several different complete datasets and combining the results.

**(b) Dealing with Noisy Data**

Noisy data refers to irrelevant data or data that has no meaning. Noise can also be referred as a random error in attribute values. Such data is generally due to collection of faulty data. The following methods can be used to deal with noisy data.

**i. Binning Method**

Binning involves dividing a continuous variable into a set of discrete intervals or bins and then converting the numerical values into categorical values representing the respective bins. The purpose of binning is to

simplify the data and reduce the noise or variations in the dataset. It can be particularly useful when dealing with large datasets or continuous variables that have a wide range of values. Binning helps to group similar values together and identify patterns or relationships that might not be apparent when using the original continuous values. There are different methods for binning in data mining, including:

1. Equal-width binning: This method divides the range of values into equal-width intervals. For example, if you have a variable ranging from 0 to 100 and want to create 5 bins, each bin would cover a range of 20 units (0-20, 20-40, 40-60, 60-80, 80-100). Let take an example: We want to create three bins: "Young," "Middle-aged," and "Elderly." The age range in the dataset is from 20 to 80 years. We can divide this range equally into three bins:

   - Bin 1: Young (20-40 years)
   - Bin 2: Middle-aged (41-60 years)
   - Bin 3: Elderly (61-80 years)

   Using this method, we have divided the continuous age variable into three discrete categories based on equal-width intervals.


2. Equal-frequency binning: In this method, each bin contains an equal number of data points. The data points are sorted in ascending order, and then the values are divided into equal-sized groups. This approach ensures that each bin has approximately the same number of observations. Here, we aim to create three bins with an equal number of individuals in each bin. Let's assume we have 100 individuals in the dataset. We can sort the ages in ascending order and divide them into three groups with approximately 33 individuals in each bin:

   - Bin 1: Young (20-35 years)
   - Bin 2: Middle-aged (36-55 years)
   - Bin 3: Elderly (56-80 years)

By dividing the ages based on equal frequencies, we ensure that each bin contains    roughly the same number of individuals.


3. Custom binning: Sometimes, it is more appropriate to define the bins based on domain knowledge or specific requirements. This method allows you to manually define the boundaries for each bin according to your understanding of the data.

Binning can be applied to both numerical and ordinal variables. However, it is important to note that binning can introduce some information loss since the original numerical values are converted into

categories. Therefore, binning should be used judiciously, taking into consideration the specific characteristics of the dataset and the goals of the data mining task.

### ii. Regression:

Regression is a statistical analysis technique used to model and examine the relationship between a dependent variable and one or more independent variables. It is commonly used for predictive analysis and understanding the influence of independent variables on the dependent variable. In regression analysis, the dependent variable is the outcome or response variable that you want to predict or explain, while the independent variables are the predictor variables that may affect the outcome. The goal of regression is to find the best-fit line or curve that minimizes the difference between the predicted values and the actual values of the dependent variable. The type of regression analysis used depends on the nature of the data and the research question. Here are a few common types of regression:

1) Simple Linear Regression: It involves a single independent variable and a linear relationship with the dependent variable. The goal is to fit a straight line to the data.
2) Multiple Linear Regression: It involves multiple independent variables and a linear relationship with the dependent variable. It aims to fit a hyperplane to the data.
3) Polynomial Regression: It models the relationship between the independent variables and the dependent variable with polynomial terms. It can capture nonlinear relationships.
4) Logistic Regression: It is used when the dependent variable is binary or categorical. It models the probability of a particular outcome using a logistic function.
   The regression methods also identify outliers and noise. The data point that do not follow the relationship pattern are outliers.

### iii. Clustering:

Clustering is a technique used in data mining and machine learning to group similar objects or data points together based on their characteristics or similarities. The goal of clustering is to discover hidden patterns or structures within a dataset without any prior knowledge of the groups or classes. In clustering, the algorithm analyzes the data and partitions it into clusters, where each cluster consists of data points that are more similar to each other compared to those in other clusters. The similarity between data points is measured using a distance or similarity metric, such as Euclidean distance or cosine similarity. The objective of clustering is to group similar data points in a cluster. Outliers or unwanted data points are easily detected in this process. The data points that do not form part of any cluster or group are outliers.

## 2. Data Transformation:

The process of data transformation is to change the data in suitable forms appropriate for data mining. The following ways can be used for data transformation:

- **Normalization**

  The process of normalization scales the data values in a specified range such as (-1.0 to 1.0 or 0.0 to 1.0). Various methods of normalization are:

  - Decimal Scaling: This method divides each value by the same power of 10. For example an attribute value ranges between -1000 and 1000, the value is changed to -1 and 1 by dividing each value by 1000.

  - Min-Max normalization: In this method the new value is computed on the basis of previously known minimum and maximum values. The new value is computed as follows-

    $$newValue = \frac{original\ value - oldMin(newMax - newMin) + newMin}{oldMax - oldMin}$$

    where *oldMax* and *oldMin* are original maximum and minimum values for the attribute under consideration. *NewMax* and *newMin* specify the new maximum and minimum values. *NewValue* represents the transformation of original value.

  - Normalization using Z-scores: This method converts a value to a standard score by subtracting the attribute mean from the value and dividing by the attribute standard deviation.

    $$newValue = \frac{original\ value - attribute\ mean\ value}{attribute\ standard\ deviation}$$

    This method is useful when minimum and maximum values of the attribute is not known.

- **Attribute Selection:**

  In this method new attributes are constructed from the already existing set of attributes.
  There are various attributes in the complete set of attributes that are highly correlated and are redundant in the set. Such attributes can be eliminated.
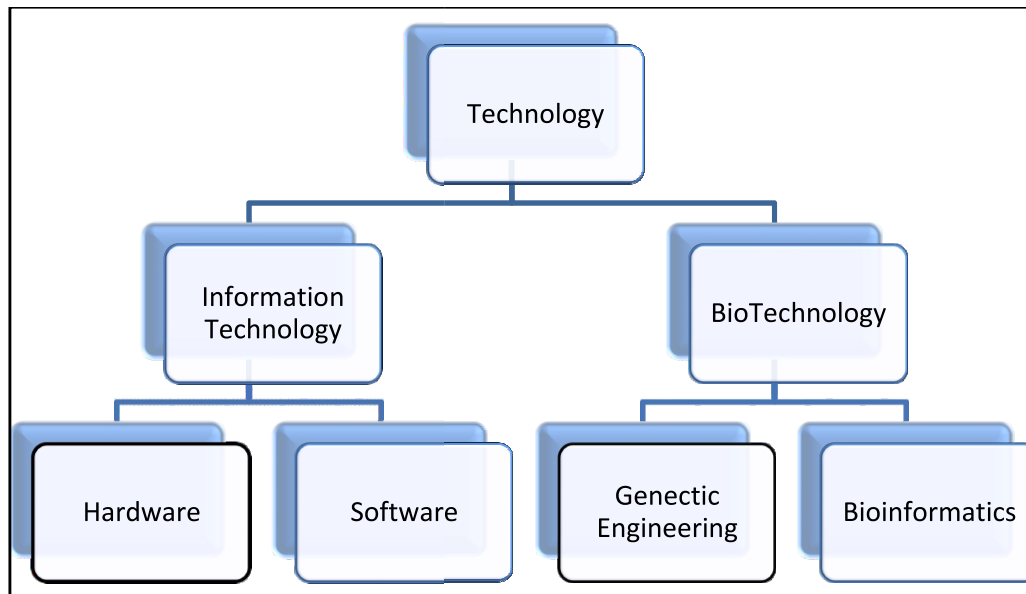
- **Discretization:**

  Discretization is the process of transforming continuous data or variables into discrete categories or intervals. This is often done to simplify the analysis and to convert data into a format suitable for certain algorithms that require discrete input. Discretization involves partitioning the range of a continuous variable into a finite number of intervals, each represented by a distinct value or category.


- **Concept Hierarchy Generation:**

  A concept hierarchy is a structured representation of knowledge in a hierarchical format, where concepts are organized from more general to more specific. It helps in understanding the relationships between different concepts and their levels of abstraction. In this process the attributes are converted from lower level to higher level in hierarchy.

  Example: Consider the concept hierarchical structure given in Figure 1.



*Figure 1: An example of Concept Hierarchy*

In the above hierarchical diagram, the main concept is "Technology." It represents the broad category of various technological fields. Underneath the main concept, there are several sub-concepts representing specific areas of technology such as "Information Technology, and " "Biotechnology," Each sub-concept is further divided into sub-concepts that provide more specific areas or applications within the respective technology field. For example, under "Information Technology," there are sub-concepts like "Hardware," "Software," and "Networks." Similarly, "Biotechnology" is subdivided into sub-concepts such as "Genetic

Engineering," and "Bioinformatics". The diagram continues this hierarchical structure, allowing for a comprehensive representation of various technology domains and their related sub-domains.

## 3. Data Reduction:

Data mining process deals with huge amount of data. Data handling is major challenge in mining process. Data reduction methods aim to overcome this challenge by reducing the data size to appropriate amount. The various ways of data reduction include:

- **Data Cube Aggregation:**

 A data cube represents data along multiple dimensions, allowing users to analyze data from different perspectives. Aggregation in a data cube involves combining and summarizing data at different levels of granularity within the dimensions. It allows for the exploration of data from various angles and the extraction of useful insights. Aggregation can be performed using different operations, such as sum, count, average, minimum, maximum, etc.

- **Attribute Subset Selection:**

  The method focuses on selecting relevant attributes from the complete set of attributes. The irrelevant and not useful attributes are discarded. The condition of relevancy can be judged by various ways. One such method is that input attributes highly correlated with other input attributes is redundant. Such attributes can be eliminated. *Example*: Consider a dataset with the following attributes:

  1. Age (numeric)
  2. Gender (categorical: male/female)
  3. Education Level (categorical: high school/college/graduate)
  4. Income (numeric)
  5. Employment Status (categorical: employed/unemployed)
  6. Marital Status (categorical: married/single/divorced)
  7. Health Condition (categorical: good/fair/poor)
  8. Loan Approval (binary: yes/no)

The goal is to predict whether a person's loan application will be approved based on their attributes. However, not all attributes may be relevant for making accurate predictions. Attribute subset selection can help identify the most informative attributes.

One approach to attribute subset selection is to use a search algorithm, such as the sequential forward selection (SFS) method. Here's how it works:

1) Start with an empty set of selected attributes.
2) Iterate through each attribute not yet selected.
3) For each attribute, train a machine learning model (e.g., logistic regression, decision tree) using the selected attributes and the current attribute under consideration.
4) Evaluate the performance of the model using a suitable evaluation metric (e.g., accuracy, area under the ROC curve) through cross-validation or a separate validation set.
5) Select the attribute that improves the model's performance the most.
6) Add the selected attribute to the set of selected attributes.
7) Repeat steps 3-6 until a predefined stopping criterion is met (e.g., a maximum number of attributes reached, no further improvement in performance).
8) Finally, train the model using the selected attributes and evaluate its performance on a separate test set.

In the example, the SFS algorithm would start with an empty set and iteratively select attributes based on their impact on the model's performance. For instance, it may find that the attributes "Age," "Income," and "Employment Status" contribute significantly to predicting loan approval, while other attributes may have minimal impact. The algorithm would then select these three attributes as the final subset for training the model.

By performing attribute subset selection, we can reduce the complexity of the model, improve interpretability, and potentially enhance the model's generalization performance by focusing on the most relevant attributes.

■ **Dimensionality Reduction:**

Data size is reduced by applying encoding methods. The reduction can be lossy or lossless. Two essential methods are Principal Component analysis (PCA) and Wavelet Transform.

PCA stands for Principal Component Analysis. It is a dimensionality reduction technique used to transform a dataset with a large number of variables into a smaller set of uncorrelated variables called principal components. These components capture the most important information or patterns in the data while minimizing the loss of information. The main idea behind PCA is to find a new coordinate system in which the data points are represented in terms of orthogonal axes called principal components. The first

principal component explains the largest amount of variance in the data, followed by the second component, and so on. Each component is a linear combination of the original variables.

The steps involved in performing PCA are as follows:

1) Standardize the data: It is important to standardize the variables to have a mean of zero and equal variances. This ensures that variables with larger scales do not dominate the analysis.

2) Compute the covariance matrix or correlation matrix: The covariance matrix represents the relationships between variables, while the correlation matrix represents the correlations. PCA can be performed using either of these matrices.

3) Compute the eigenvectors and eigen values: The eigenvectors represent the directions or axes of the principal components, and the eigen values represent the amount of variance explained by each component. They are obtained by solving the eigen value problem associated with the covariance or correlation matrix.

4) Sort the eigen values and select the top-k components: The eigen values are sorted in descending order, and the corresponding eigenvectors are selected to form the principal components. The number of components to retain depends on the desired level of dimensionality reduction.

5) Transform the data: The original data is projected onto the new coordinate system defined by the selected principal components. This transformation yields a reduced-dimensional representation of the data.

PCA has several applications, including data visualization, feature extraction, and noise reduction. It helps in identifying the most influential variables, detecting patterns or clusters in the data, and improving computational efficiency by reducing the dimensionality of the data.

---

*Check your progress 2*

*1.What is the significance of data-preprocessing in data mining?*

*2.What is data cleaning?*

*3.Write names of different methods to deal with noisy data.*

*4. Differentiate equal-width binning and equal-frequency binning methods.*

*5. How is the regression technique helpful to reduce the noise?*

## 1.4 Measures of Similarity

Similarity measures are numerical measures of the degree of similarity or alikeness between objects. The measure of similarity is known as proximity. The objective of similarity measures is to compute distance between data points. The distance measure is a measure of similarity of similarity or dissimilarity between data. The smaller the distance between the objects, the more is the similarity between the objects. Various measures such as Euclidean distance and Manhattan are measures of similarity.

Similarity measures in data mining are used to quantify the similarity or dissimilarity between objects or data instances. These measures play a crucial role in various data mining tasks, such as clustering, classification, recommendation systems, and information retrieval. Here are some commonly used similarity measures in data mining:

➢ Euclidean Distance: Euclidean distance is a popular similarity measure used for numeric data. It calculates the straight-line distance between two points in a multidimensional space. For two data instances represented as vectors, the Euclidean distance is computed as the square root of the sum of squared differences between corresponding attributes.

➢ Cosine Similarity: Cosine similarity is often used for text and document analysis. It measures the cosine of the angle between two vectors, which represents the similarity of their orientations in the feature space. It is particularly useful when the magnitude of the vectors is not important, and the focus is on the direction or orientation of the vectors.

➢ Jaccard Similarity: Jaccard similarity is a measure commonly used for binary or categorical data. It calculates the ratio of the intersection of two sets to the union of the sets. It is particularly useful for measuring the similarity between sets, such as in collaborative filtering or item-based recommendation systems.

➢ Hamming Distance: Hamming distance is a similarity measure used for categorical or binary data. It calculates the number of positions at which two strings of equal length differ. It is often used in data mining tasks involving sequences, such as DNA analysis or text mining.
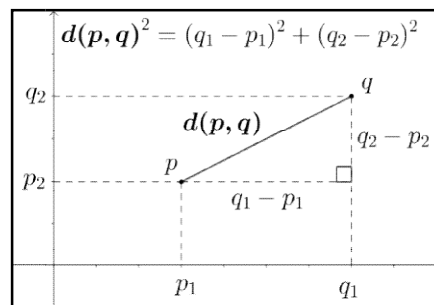
➢ Manhattan Distance: Manhattan distance, also known as city block distance or L1 norm, is a distance measure used for numeric data. It calculates the sum of absolute differences between the coordinates of two points. It is particularly suitable when movement is constrained to a grid-like structure or when attribute values are measured in different units or scales.

➢ Minkowski Distance: Minkowski distance is a generalized distance measure that includes both Euclidean distance and Manhattan distance as special cases. It is defined as the p-th root of the sum of the p-th powers of the differences between corresponding attributes. By varying the value of the parameter p, different types of Minkowski distances can be obtained.

➢ Pearson Correlation Coefficient: Pearson correlation coefficient is a measure of the linear relationship between two variables. It ranges from -1 to 1, where -1 indicates a perfect negative correlation, 1 indicates a perfect positive correlation, and 0 indicates no correlation. It is commonly used for assessing the similarity between numeric data instances.

These are some of the widely used similarity measures in data mining. The choice of the similarity measure depends on the characteristics of the data and the specific requirements of the data mining task at hand.

1.4.1 Euclidean Distance

Euclidean distance computes the root of squared differences between the coordinates between two objects. In the Euclidean plane, let point p have Cartesian coordinates (p1,p2) and let point q have coordinates (q1,q2). Then the distance between p and q is given by applying the Pythagorean theorem to a right triangle with horizontal and vertical sides, having the line segment from p to q as its hypotenuse. Thus Euclidean distance is defined as follows (Figure 2):



*Figure 2: Euclidean distance*

Source: https://en.wikipedia.org/wiki/Euclidean_distance#/media/File:Euclidean_distance_2d.svg

## 1.4.2 Manhattan Distance

The Manhattan distance is another way of measuring distance. Manhattan distance, also known as city block distance , is a distance metric used to measure the absolute difference between two points in a grid-like structure. It is named after the Manhattan grid system, where distances are calculated by following the perpendicular city blocks.

If there are two points, (x1, y1) and (x2, y2) the Manhattan distance is defined as-

$$|x1 - x2| + |y1 - y2|$$

For example, let's consider two points in a 2D space: A(2, 3) and B(5, 7). To calculate the Manhattan distance between these two points, we apply the formula. Manhattan distance = |2 - 5| + |3 - 7| = 3 + 4 = 7. Thus the Manhattan distance between points A (2, 3) and B(5, 7) is 7.

## 1.4.3 Similarity and Dissimilarity matrix

A similarity matrix and a dissimilarity matrix are used to quantify the similarity or dissimilarity between pairs of data objects. These matrices are commonly employed in various data mining tasks, such as clustering, classification, and recommendation systems. A similarity matrix measures the similarity between objects, where higher values indicate greater similarity. On the other hand, a dissimilarity matrix quantifies the dissimilarity between objects, where higher values indicate greater dissimilarity or distance. The choice between using a similarity or dissimilarity matrix depends on the specific data mining task and the nature of the data being analyzed.

 A similarity matrix is a square matrix where each element represents the similarity between a pair of objects. It is symmetric, as the similarity between object A and object B is the same as the similarity between object B and object A. Similarity measures can vary depending on the data and the context of the problem. Common similarity measures include cosine similarity, Euclidean distance, Jaccard similarity, and Pearson correlation coefficient. The values in the similarity matrix typically range from 0 to 1, with 1 indicating maximum similarity.  A Dissimilarity Matrix: A dissimilarity matrix, also known as a distance matrix, quantifies the dissimilarity or distance between pairs of objects. Unlike a similarity matrix, a dissimilarity matrix need not be symmetric. Dissimilarity measures are typically based on metrics such as Euclidean distance, Manhattan distance, Minkowski distance, or Hamming distance, depending on the nature of the data. The values in the dissimilarity matrix can vary depending on the specific distance measure being used.

Both similarity and dissimilarity matrices are used in different data mining algorithms and techniques:

➢ Clustering algorithms, such as k-means or hierarchical clustering, often utilize dissimilarity matrices to group similar objects together based on their distances or dissimilarities.

➢ Classification algorithms, such as k-nearest neighbours (KNN), may employ similarity matrices to find the most similar instances or neighbours for making predictions.

➢ Collaborative filtering techniques in recommendation systems may use similarity matrices to identify similar users or items based on their preferences or behaviour.

➢ Dimensionality reduction techniques like multidimensional scaling (MDS) or principal component analysis (PCA) can use either similarity or dissimilarity matrices to project high-dimensional data into lower-dimensional spaces while preserving the relationships between objects.

Overall, both similarity and dissimilarity matrices play a crucial role in data mining tasks, enabling the analysis of patterns, relationships, and similarities/dissimilarities among data objects, ultimately assisting in extracting valuable insights from the data.

Example of similarity and dissimilarity matrix-

Let's consider a similarity matrix example for a dataset of objects based on their pairwise similarities. In this case, we'll use a simple binary similarity measure, where 1 indicates a similarity and 0 indicates dissimilarity. Suppose a dataset of four objects: A, B, C, and D. We want to measure the similarity between these objects based on a specific criterion, such as their color. Here's an example similarity matrix based on color similarity:

A 1 0 1 1

B 0 1 0 0

C 1 0 1 1

D 1 0 1 1

In this example, we have compared the objects' color similarity with a binary measure. The values in the matrix indicate whether two objects have similar (1) or dissimilar (0) colors. From the similarity matrix, we can observe the following:

▪ Object A is similar to objects C, D, and itself.

- Object B is similar only to itself.
- Object C is similar to objects A, C, and D.
- Object D is similar to objects A, C, and D.

The values in the similarity matrix are based on the specific similarity measure used and the data provided in this example.
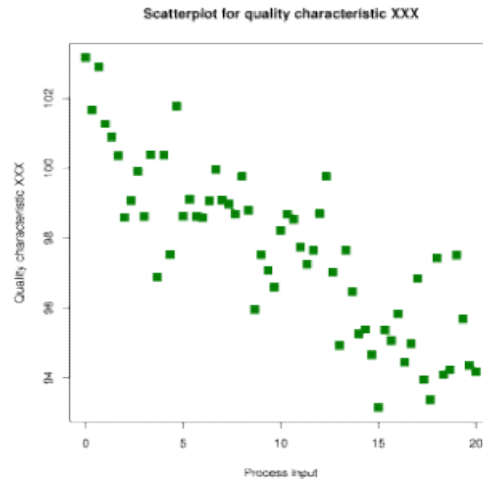
# 1.5 Data Visualization

Data visualization is the graphical representation of data and information using visual elements such as charts, graphs and maps. Its purpose is to present complex data sets in a visual format that is easier to understand, interpret, and derive insights from. Data visualization is an essential part of data analysis and communication in various fields, including business, science, research, finance, and journalism. It allows individuals to explore patterns, trends, and relationships in data, uncover hidden insights, and effectively communicate their findings to others. The objective of data visualization is to represent data with graphics, charts, and plots. The information is represented as patterns of data and various outliers. Data visualization is the implication of phrase "A picture is worth a thousand words". Various complex relationships between data and data insights can be easily understood and analyzed by a visual representation such as a graph. Various useful patterns and data trends can be easily identified by a visual graph or plot. In text mining a word cloud representation easily demonstrates essential concepts and hidden relationships within unstructured data. A knowledge graph can be used to represent relationships between entities.

Dashboards are data visualization tools for visualizing data across multiple data sources. Some visualization techniques are:

- **Tables:** A table comprises of rows and columns used to compare variables. Tables depict information in a structured way.
- **Bar charts:** These charts use rectangular bars to represent data values. They are useful for comparing data across different categories or time periods.
- **Pie charts:** These graphs are separated into sections. Each section represents a part of a whole. They are ideal for showing proportions or percentages.
- **Line charts:** A line chart plots data points over time and are useful for predictive analysis.
- **Histograms:** The graph represents a distribution of numbers using a bar chart (with no spaces between the bars).
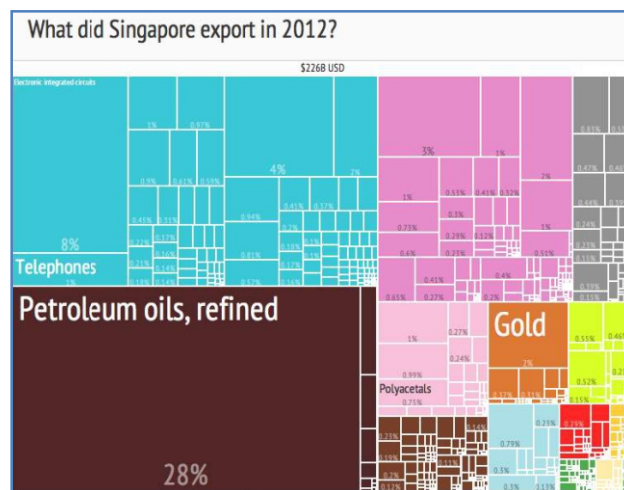
- **Scatter plots:** The graph plots two variables of data along the *x* and *y* axes. It is used to show the relationship between two quantitative variables. Scatter plots use a two-dimensional grid to display individual data points as dots. They are helpful for visualizing the relationship between two variables and identifying patterns or correlations. An example scatter plot is given in Figure 3.



*Figure3: Scatter Plot*

*Source: https://en.wikipedia.org/wiki/Scatter_plot*

- **Tree maps** display hierarchical data as a set of nested shapes. Tree maps are useful for comparing the proportions between categories via their area size. An example tree map is given in figure 4.



*Figure 4: Tree map showing Singapore exports in 2012*
*Source: https://en.wikipedia.org/wiki/Treemapping*

To create data visualizations, some specialized software and tools are Tableau, Microsoft Power BI, Python libraries like Matplotlib and Seaborn, R programming language with ggplot2, and web-based tools like D3.js. Effective data visualization requires careful consideration of the audience, the purpose of the

visualization, and the type of data being presented. It's important to choose the most appropriate visualization type that effectively represents the data and facilitates understanding and insight generation.

1. What is data visualization?

2. Write names of some data visualization technique.

3. Write names of Python libraries that provide data visualization support.

4. What is histogram?

5. How is the scatter plot useful?

## 1.6 Summary

In this unit basic concepts of data mining were presented. Representation of data according to its nature is important for analysis and mining tasks. There are various data types that can be used to represent data according to specific requirements. The essential data types include nominal, ordinal, discrete, and continuous data. Usually, a raw data needs preprocessing before the mining process. The data pre-processing is the preliminary phase of the mining process that is useful to improve data quality and transform data in a suitable format. Data cleaning is an important pre-processing task that is required to reduce the noise and other anomalies in the data. Data transformation techniques are applied to make the data appropriate for the mining. Major transformation techniques include data normalization, discretization, concept hierarchy generation, and attribute selection, etc. Data mining deals with a huge amount data. Sometimes, it is desirable to reduce the data dimensions for the better performance. For this purpose, dimensionality reduction techniques are applied that reduce the data dimensions while retaining the important information. PCA and wavelet transform are the major dimensionality reduction techniques. Similarity measures play a crucial role in data mining to quantify the similarity or dissimilarity between data points or objects. Different types of distance measures including Euclidian distance, Manhattan distance, and Hamming distance, etc. can be used as similarity measures. Data visualization is another important aspect of data mining that allows to analyse the data with help of graphs and charts. There are specialized tools and libraries available that facilitate to effectively represent data for understanding and insight generation.

In this unit we presented various data types used in Data Mining process. We also described the need of data-preprocessing and data visualization. Various methods used in data-pre-processing and data

visualization were also discussed in this block. Section 1.4 discussed various measures of similarity. To summarise the unit focuses on the following:

- ✓ Essential types of data used in Data mining process
- ✓ The purpose of data preprocessing/data preparation in data mining.
- ✓ Various methods used for data-preprocessing
- ✓ The importance of Data Visualization in Data mining.
- ✓ Various methods and diagrams for Data visualization
- ✓ The concept of similarity and dissimilarity in data mining.
- ✓ Various measures of similarity to measure distance between objects

---

## 1.7 REVIEW QUESTIONS

---

Q1. What is Data Mining? What is the relation between data warehousing and data mining.

_____

_____

Q2. Explain the differences between Knowledge discovery and data mining.

_____

_____

Q3. What do you mean by Data Pre-processing? Discuss various methods of Data Pre-processing.

_____

_____

Q4. Explain the importance of data cleaning in data mining process. Describe different data cleaning approaches. How can we handle missing values?

_____

Q5. Explain Noisy Data. Discuss methods to handle noisy data in Data preprocessing phase.

_____

_____

Q6. Explain the significance of similarity measures in Data Mining. Discuss at least three similarity measures.

_____

_____

Q7. Discuss data imputation as an essential data preprocessing step.

_____

_____

Q8. Describe the different types of data that can be used in data mining. How do you handle each type?

_____

_____

Q9. What is clustering? What is the significance of similarity measures in clustering?

_____

_____

Q10. Explain at least 5 ways/methods for data visualization.

_____

_____

# UNIT 2 : **Data warehousing - I**

**Structure**

## 2.0 Introduction

In previous unit of this block, you learned about essential data types used in data mining process. The various data-preprocessing methods used in data mining process were also discussed. Similarity measures and data visualization were also discussed in the previous unit. This unit will focus on the concept of data warehousing, various components of data warehousing and integration with data mining process. Data warehouse is valuable repository for business organisations for smart and effective decision making. Data warehouse implementation and various topics related to data warehousing is the covered in this unit.

The need of data warehouse and the various characteristics of a data warehouse are explained in this unit. The concept of data marts and how they are different from a data warehouse is explained in this unit. The client server computing model and it's implementation for a data warehouse is explained in this unit.

## 2.1 Objectives

This unit deals with the representation of algorithm in terms of flowchart and pseudo code.

At the end of this unit, you will be able to:

- Explain the need of data warehousing

- Understand the characteristics of a data warehouse

- Understand Data Marts and differentiate between data mart and data warehouse

- Understand the Client  Server computing model and data warehouse

## 2.2 Data Warehousing

In words of Bill Inmon,*" A warehouse is a subject-oriented, integrated, time-variant and non-volatile collection of data in support of management's decision making process"*.

Data warehouse is the integration of data from multiple heterogeneous sources of data that integrate business data for efficient decision making. A data warehouse is a centralized repository that stores large amounts of structured and/or unstructured data from various sources within an organization. It is designed to support business intelligence, reporting, and data analysis activities. The data warehouse integrates data from different operational systems, such as transactional databases, spreadsheets, and external sources, to provide a unified view of the organization's data.

A data warehouse is different from databases and operational data sources that focus on day to day operations of the organization and not on decision making and strategic planning.

Data warehouse contains historical data collected over past years for analysis of business data. The essential characteristics of a data warehouse are:

- *Subject Oriented*-A data warehouse provides information related to themes or subjects such as Sales, Product, and Customer of business. It does not provide information related to operational functionalities of business organizations.
- *Integrated*- A data warehouse is a repository of data collected from various heterogeneous data sources. A data warehouse consolidates data from multiple sources into a unified and consistent format. This integration process involves transforming and cleaning the data to ensure its quality and consistency.
- *Time Variant*- The data in the data warehouse is not current data. The data is loaded into the warehouse through an Extract, Transform, Load (ETL) process, and once stored, it remains unchanged. This ensures data consistency and provides a stable environment for reporting and analysis. It is repository of historic data as the objective of data warehousing is analysis and decision making.
- *Non-Volatile*- The data in data warehouse is added to previous existing data. No previous data is removed or changed. The emphasis is on enhancing the amount of data over the years for decision making.

The essential advantages of using a data warehouse include:

➢ Centralized Data: Data warehouses provide a single, consolidated view of data from various sources, eliminating data silos and inconsistencies.

➢ Improved Data Quality: Data integration processes in a data warehouse help to standardize and cleanse data, ensuring better data quality and accuracy.

➢ Enhanced Reporting and Analysis: Data warehouses enable easy retrieval and analysis of data for reporting, ad-hoc queries, and business intelligence purposes. They support complex queries and provide faster response times.

➢ Historical Analysis: By storing historical data, data warehouses enable trend analysis, forecasting, and comparison of data over time.

➢ Decision Making: Access to accurate and comprehensive data in a data warehouse empowers decision-makers to make informed decisions based on reliable insights.

Examples of popular data warehouse technologies include

• IBM DB2 Warehouse,
• Microsoft Azure Synapse Analytics,
• Amazon Redshift, and Snowflake.

These technologies provide robust platforms for storing and analyzing large volumes of data in a data warehouse environment.

*Check your Progress 1*

*1. What is a Data Warehouse and why is it important?*
*2. Explain the essential characteristics of a Data warehouse.*
*3. What are advantages of data warehouse?*
*4. Give some examples of data warehouse technologies.*

## 2.3 Integration with Data Mining

The data warehouse functions as a huge data repository for data mining process. The ability of a data warehouse to handle huge and historical amount of data makes it a essential source of data in data mining process for extracting relevant, useful information for decision making.

Data Warehouses can easily handle today's highly scalable business and increased volumes of day-to-day transaction data required for better and efficient decision making. A data warehouse provides access to

historical data collected over past many years for finding efficient patterns of information. A data warehouse is not restricted to internal data repositories. It is expandable on cloud platforms. Thus data is available from vast resources leading to better productivity. Data being available in a central repository allows availability from the central repository and faster decision making.

The Data warehouse provides a centralized data source for data mining process. A Data mining process model is four step process comprising of the following steps:

- Collecting data form Data Warehouse

- Presenting the data to Data mining tool/software

- Interpretation of results

- Analysis and application of results

The data mining method is applied to the data in the data warehouse to generate results and perform decision making on the generated results. An example is Market Basket Analysis. The historical data of supermarkets is maintained in data warehouse. Associations Rule Mining is applied to data to generate association rules between data items. Thus integrating data warehouse as a data repository is useful for mining results from data.
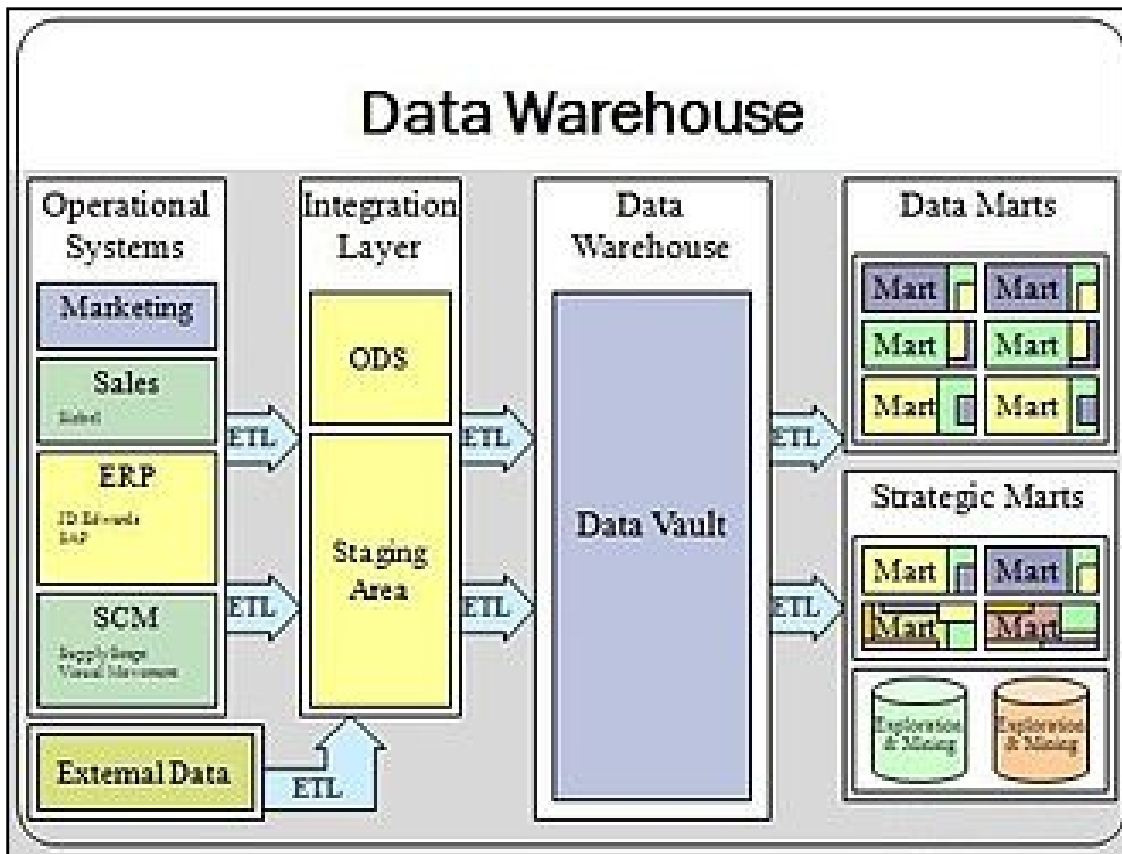
## 2.4 Data Marts

A Data mart is a subset of a Data Warehouse. It is focused on a single subject or single department. This contrasts with a Data Warehouse in which data is available from multiple subjects. A Data Mart covers the requirements of a particular business unit such as finance or sales department. A data mart enhances user responses and also decreases the volume of data for data analysis. A data mart is a subset of a data warehouse that focuses on a specific subject area or department within an organization. It is designed to provide a more targeted and streamlined view of data for specific business functions or user groups. Unlike a data warehouse, which consolidates data from various sources across the organization, a data mart typically contains data that is relevant to a particular department or business unit. The essential characteristics of a data mart include:

1. *Subject-Specific:* A data mart is centered around a specific subject area, such as sales, marketing, finance, or human resources. It focuses on providing data and analytics tailored to the needs of a specific department or user group.

2. *Subset of Data Warehouse:* A data mart is derived from a data warehouse but contains a subset of the data. It selects and organizes relevant data from the data warehouse to provide a more focused and efficient view.

3. *Simplified Data Structure:* Data marts are designed with a simplified and optimized data structure that aligns with the specific reporting and analysis requirements of the subject area. This structure facilitates faster query performance and easier access to relevant information.

4. *Departmental Ownership:* Data marts are often developed and maintained by individual departments or business units. They have ownership over the data within their specific data mart, allowing them to customize and manage the data to meet their specific needs.

5. *Increased User Accessibility:* Data marts provide a user-friendly interface for accessing and analyzing data. They are designed with the specific needs and skill sets of the users in mind, making it easier for them to perform data analysis and generate insights.

*Benefits of using data marts include:*

- Improved Performance: By focusing on a specific subject area, data marts provide faster query response times and analysis for the targeted department or user group.

- Increased Relevance: Data marts deliver information that is directly relevant to the specific business functions or user needs, making it easier to extract actionable insights.

- Simplified Data Access: Users can access and analyze data within their data mart without needing to navigate through the entire data warehouse, improving efficiency and usability.

- Enhanced Data Security: Data marts offer tighter control over data access, as each department or user group can define and manage their own access rights and security measures.

- Agility and Flexibility: Data marts allow departments to have greater control and autonomy over their data, enabling them to adapt quickly to changing business requirements and analysis needs.

*Figure 1:Data Marts as subset of Data-warehouse Environment*
Source: https://en.wikipedia.org/wiki/Data_mart

2.4.1 Data Marts vs. Data Warehouse

A data mart is different from a data warehouse. The following are differences between them.

- ❖ Data mart is for a particular company unit or department. It is a subset of a data warehouse.

- ❖ Data marts handle department or unit level tactical decisions whereas a data warehouse handles strategic decisions of the organization.

- ❖ Data warehouse is a repository of multiple data sets and takes time to update. On the other hand data marts handle smaller data sets.

- ❖ The implementation of a Data warehouse takes many years. Data marts are smaller in scope and implemented in months.

2.4.2 Advantages of Data Marts

A data mart has certain merits and they are as follows:

❖ **Easy and fast access to data:** Data marts contain subset of data of a data warehouse. This provides easy and fast access to required data.

❖ **Fast decision making:** Data marts are focused on a particular subject or unit of the organization. This allows faster analysis and decision making.

❖ **Lower cost.** The cost of setting a data mart is much lower as compared to a Data Warehouse.

❖ **Easy implementation & maintenance.** The implementation and maintenance of a data mart is relatively easy as compared to a Data Warehouse.

❖ **Better data access control.** Data in data mart is subset of data warehouse. This allows controlling data access privileges at a granular level.

---

*Check your progress 2*

*1. What is a data mart?*
*2. How is data mart different from Data Warehouse?*
*3. What are the essential characteristics and advantages of Data Mart?*
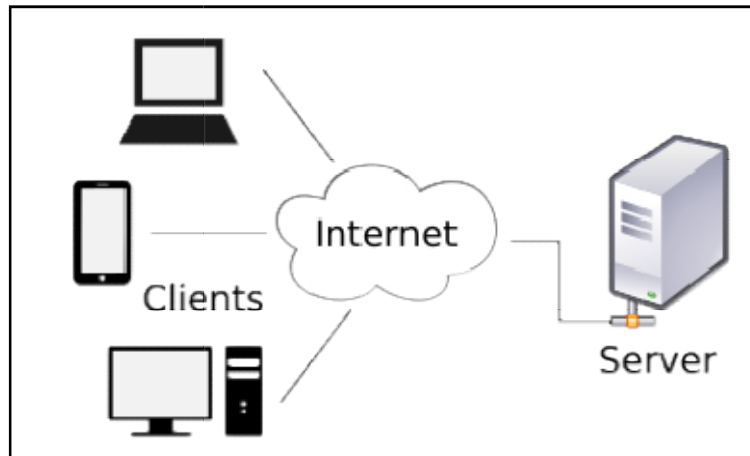*4. What are advantages of Data Mart?*

---

## 2.5 Client Server Computing Model and Data Warehousing

The client server computing model is a distributed architecture in which the workload is distributed between servers and clients. The clients do not share their resources. The server resources are shared among the clients. The clients send requests to the servers. The servers process the requests of the clients and reply to the clients.

In client-server architecture, the system is divided into two main components: the client and the server. The client component represents the user interface and application running on the user's device, while the server component hosts the application logic and data.
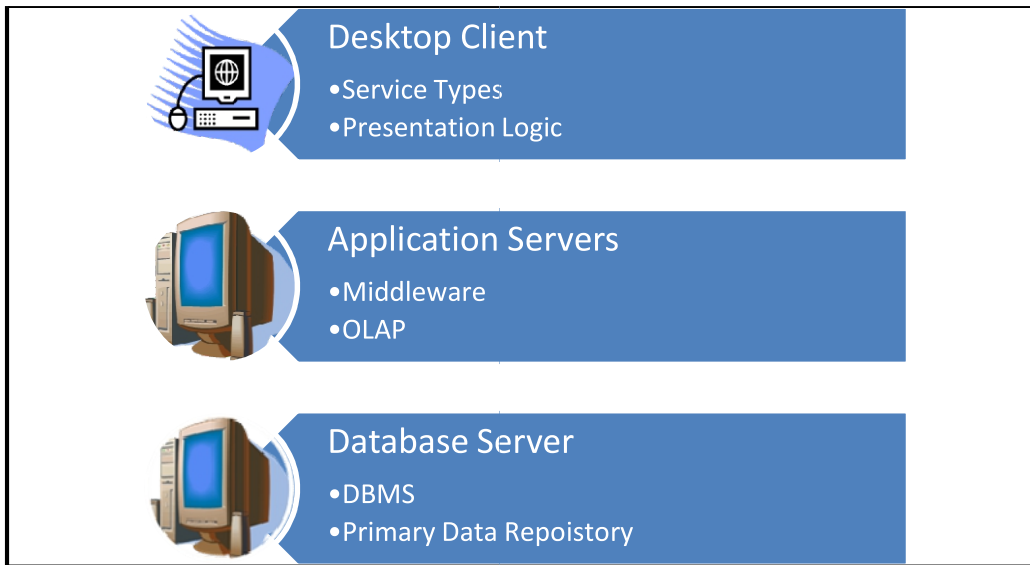
*Figure 2: Client-Server Architecture*

The interaction between the client and the server takes place as follows:

- The client component includes the user interface, which is responsible for presenting information and interacting with the user. It can be a desktop application, a web browser, or a mobile app.

- The client application communicates with the server-side application by sending requests and receiving responses. These requests can include data retrieval, updates, or other operations.

- On the server side, there is the application server, which handles the business logic and processes the client's requests. It performs the necessary computations, accesses data from the database, and generates responses to be sent back to the client.

- The server-side application may interact with a database server to store and retrieve data. The database server manages the storage and retrieval of data requested by the server-side application. The client server computing model functions with protocol of request and reply. The client sends a request to the server. The server responds with the desired information. An example of a client server computing system is a web server. It returns the web pages to the clients that requested them. Data warehousing in today's time is implemented using Client/Server Architecture. The implementation of a data warehouse is multi-tiered, second generation client-server architecture. The following diagram represents Client-Server architecture for the Data Warehouse-

*Figure 2: Client Server architecture for Data Warehouse implementation*

In the above diagram the Desktop client interacts with the Application Servers. The application servers perform analytical processing. The Database servers are the primary data repository for the Data warehouse environment.

---

*Check your Progress 3*

*1. What are the essential features of client-server architecture?*

*2. Discuss the implementation of Data warehouse as a client server architecture.*

---

## 2.6 Summary

The previous unit focused on the essential and basic concepts of data types, data pre-processing, and similarity measures and data visualization in Data mining domain. In this unit, we studied about Data warehousing concepts and integration of Data warehouse with Data Mining. The unit focused on the following:

- Description of a Data Warehouse

- Characteristics of A Data warehouse

- Advantages of a Data Warehouse

- Integration of Data warehouse and Data Mining

- Description of Data Marts

- Client-server architecture and implementation of Data warehouse as Client-server architecture

## 2.7 Review Questions

Q1. Explain the concept of a data warehouse.

_____

_____

Q2. What is a data mart? Explain its functionality and implementation.

_____

_____

Q3. Discuss the characteristics of a data warehouse.

_____

_____

Q4. Explain the various advantages of a data warehouse.

_____

_____

Q5. Discuss the implementation of a data warehouse as client server architecture    .

_____

_____

# UNIT 3: Data warehousing-II

**Structure**

## 3.0 Introduction

The unit enhances the description of Data warehouse with emphasis on essential components of a Data warehouse and layered architecture of Data warehouse. In this unit, we will learn about the essential modules of load manager, warehouse manager, Query manager and end-user tools. The 3-layered architecture of a Data Warehouse comprising of database, analytical processing module and front-end tools is also discussed in this unit. Various implementation issues related to implementation of Data warehouse is also discussed in this unit. The hardware and software requirements for implementing Data warehouse is also covered in this unit. The unit lays emphasis on data warehouse. The definition of a data warehouse, its various components, the layered architecture of a data warehouse are explained in details.

## 3.1 Objectives

As discussed in the previous unit a Data Warehouse is a centralized data repository of multiple heterogeneous data sources. It is the integration of various data sources that provides a holistic view of the organization information. It enables efficient mining of useful, valuable, previously unknown, hidden patterns of information for decision making. At the end of this unit, you will be able to:
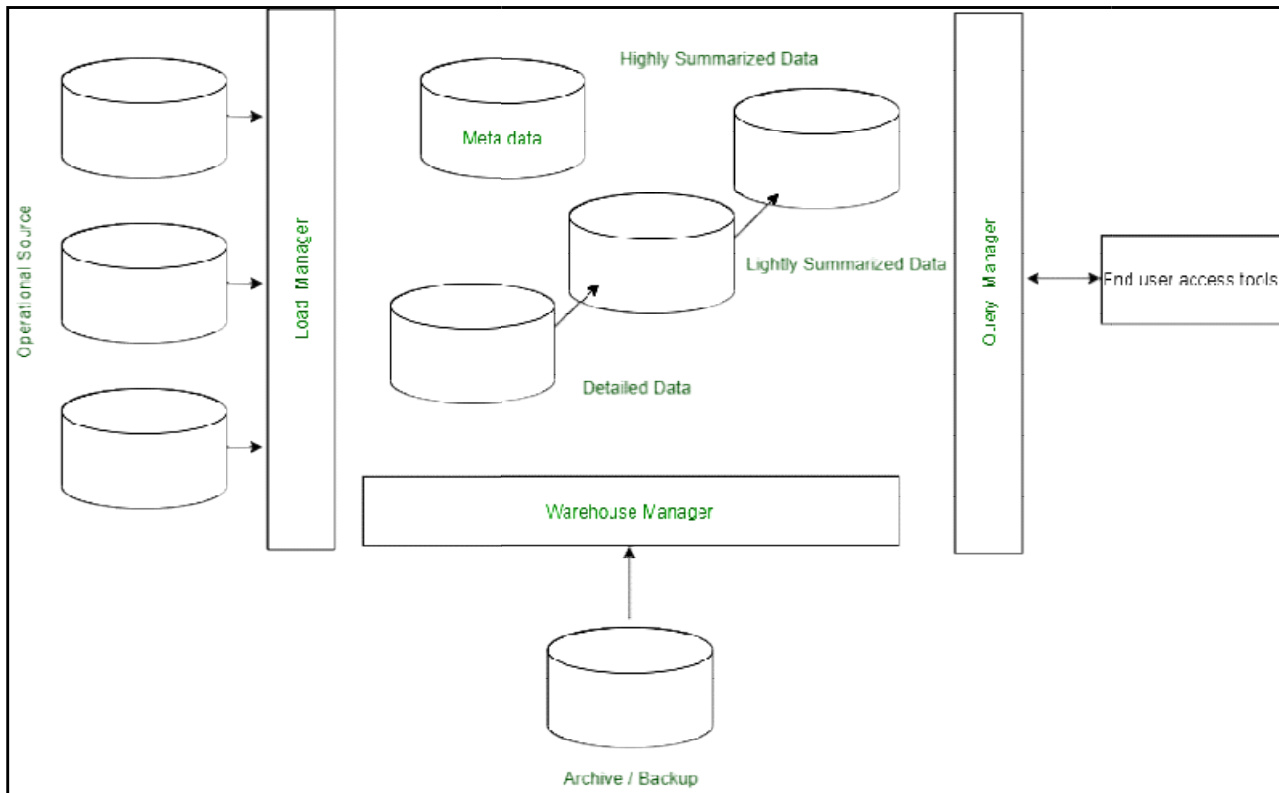
- Define a Data Warehouse

- Explain the essential components of a Data warehouse

- Explain the 3-layered architecture of Data warehouse

- Discuss the implementation issues related to Data warehouse

## 3.2 Components of a Data Warehouse

A data warehouse is not a single entity. It is composed of various components and each component has a specific signifance in the data warehouse. The essential components of a Data Warehouse are represented in the following figure-



*Figure 1: Data Warehouse components*

**Source: https://www.geeksforgeeks.org/implementation-and-components-in-data-warehouse/**

- Load Manager – This component collects data from operational systems and changes into the form useful for users. It is related to import and export of data from operational systems.
- Warehouse Manager—This is the data warehouse itself. It is the main component that contains the data from multiple heterogeneous data sources.
- Query Manager— This component provides end-users access to the stored information in data warehouse via various tools. Various access tools include OLAP (Online Analytical Processing Tools), graphical tools, and geographical information systems.
- End User Tools— These tools belong to various categories such as Query tools, Data reporting tools, OLAP tools etc.

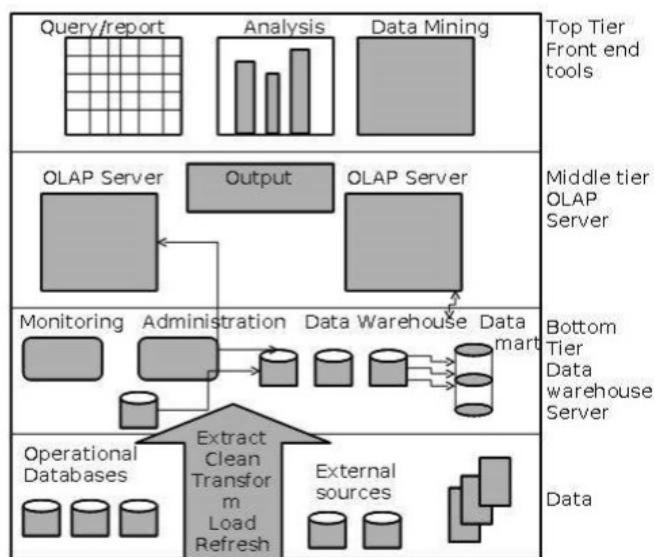## 3.3 Three-Level Architecture of Data Warehouse

Data warehouse architecture defines design framework for data storage. The data warehouse architecture plays a very important function in the data enterprise. The objective of data warehousing is to collect data from multiple heterogeneous sources and convert it into a unified form. The 3-level architecture of a data warehouse is a conceptual framework that describes the organization and structure of data within a data warehousing system. It consists of three main levels:

- the bottommost level is the data source or operational database,
- the middle level is the data warehouse itself,
- the topmost level is the data presentation layer.

➤ Operational Database/Source Systems: At the lowest level, you have the operational database or source systems. These are the systems that capture and store the transactional data generated by an organization's day-to-day operations. Examples of source systems include customer relationship management (CRM) systems, enterprise resource planning (ERP) systems, and transactional databases. The operational database is designed for efficient transaction processing and typically optimized for online transaction processing (OLTP) workloads.

➤ Data Warehouse: The middle level is the data warehouse, which serves as a centralized repository for storing and organizing data from various source systems. The data warehouse is designed for efficient querying and analysis, making it suitable for online analytical processing (OLAP) workloads. It undergoes a process known as extract, transform, load (ETL), where data is extracted from the source systems, transformed into a consistent and coherent format, and then loaded into the data warehouse.

The data warehouse typically consists of the following components:

- Staging Area: It is an intermediate storage area where data from the source systems is initially loaded before further transformation and integration into the data warehouse.
- Data Integration Layer: This layer performs data cleansing, aggregation, consolidation, and other transformations to ensure data quality and consistency.
- Data Storage Layer: It comprises one or more databases that store the integrated and transformed data in a structured format optimized for analytical queries.
- Metadata Repository: It stores information about the data in the data warehouse, including data definitions, relationships, and business rules.
- Indexes and Aggregates: These are created to improve query performance by pre-calculating and summarizing data.

➢ Data Presentation Layer: The topmost level is the data presentation layer, which provides an interface for users to access and analyze the data stored in the data warehouse. It includes various tools and technologies for reporting, visualization, and analysis, such as dashboards, ad hoc query tools, and data mining tools. The data in this layer is presented in a format that is easy to understand and supports decision-making processes within the organization. The data presentation layer enables users to perform tasks such as generating reports, creating visualizations, conducting ad hoc queries, and performing advanced analytics on the data stored in the data warehouse.



*Figure 2: 3-layered architecture of Data Warehouse*
*Source: https://www.tutorialspoint.com/dwh/dwh_architecture.htm*

In the above diagram the three layers are as follows-

▪ **Bottom Tier/Layer:**

It is the database where data is loaded.

▪ **Middle tier/Layer:**

It is the application layer that provides abstract view of the system. OLAP server is used for analysis. OLAP server is implemented using ROALP and MOALP server.

An OLAP (Online Analytical Processing) server is a high-performance database system optimized for complex queries and analytical tasks. OLAP servers are designed to support multidimensional analysis of

large datasets, making them ideal for business intelligence (BI) applications, data mining, and decision support systems. The various types of OLAP Servers are,

1. **MOLAP (Multidimensional OLAP)**

   - **Storage:** Data is stored in a multidimensional cube.
   - **Performance:** Fast query performance due to pre-aggregated data.
   - **Example:** Microsoft Analysis Services, Oracle Essbase.

2. **ROLAP (Relational OLAP)**
   - **Storage:** Data is stored in relational databases.
   - **Performance:** Slower than MOLAP but can handle large datasets.
   - **Example:** IBM Cognos, MicroStrategy.

3. **HOLAP (Hybrid OLAP)**
   - **Storage:** Combines MOLAP and ROLAP approaches.
   - **Performance:** Balances the benefits of both MOLAP and ROLAP.
   - **Example:** Microsoft Analysis Services (can operate in HOLAP mode)

- **Top tier/Layer:**
   It is the user interaction layer. It is the front-end that uses reporting tools, query and analysis tools.

The 3-level architecture of a data warehouse provides a structured and organized approach to store, integrate, and present data for decision support and analytical purposes, separating the operational systems from the analytical processes to ensure optimal performance and flexibility.

> *Check your progress 1*
>
> *1. Explain the essential components of a Data warehouse.*
>
> *2. Explain the 3-layered architecture of a Data warehouse with a diagram.*

## 3.4 Implementation Issues In Data Warehousing

Data warehousing is a complex process that involves the collection, integration, transformation, and storage of data from various sources to support decision-making and analysis. While data warehousing offers numerous benefits, there are several implementation issues that organizations may encounter. Here are some common challenges in data warehousing implementation:

❖ Data quality and consistency: Data from different source systems may have inconsistencies, errors, or missing values. Ensuring data quality and consistency during the integration process is crucial for accurate reporting and analysis. Data cleansing and transformation techniques, such as data profiling, standardization, and validation, need to be implemented to address these issues.

❖ Data integration: Organizations often have multiple source systems with different data formats, structures, and naming conventions. Integrating data from disparate sources into a unified format can be challenging. Extract, Transform, Load (ETL) processes are commonly used to extract data from source systems, transform it to match the data warehouse schema, and load it into the data warehouse.

❖ Scalability and performance: As the amount of data grows, the performance of the data warehouse can be affected. Query response times may increase, leading to slower reporting and analysis. Implementing proper indexing, partitioning, and optimizing the data warehouse schema can help improve performance. Additionally, organizations may need to consider scaling options such as distributed processing or cloud-based solutions to handle large volumes of data.

❖ Data security and privacy: Data warehousing involves storing and accessing sensitive and valuable data. Ensuring data security and privacy is essential to protect against unauthorized access, data breaches, and compliance violations. Organizations need to implement appropriate security measures, including user authentication, data encryption, access controls, and audit trails, to safeguard the data warehouse.

❖ Change management and governance: Data warehousing initiatives often involve significant changes to data management processes and organizational workflows. Managing these changes and ensuring stakeholder buy-in can be challenging. Establishing clear governance policies, defining roles and responsibilities, and providing training and support for end-users are crucial for successful implementation.

❖ Data governance and metadata management: Data governance encompasses the policies, processes, and standards for managing and ensuring the quality, availability, integrity, and security of data. Establishing effective data governance practices, including metadata management, data stewardship, and data lineage tracking, is essential for maintaining data integrity and supporting data-driven decision-making.

❖ Cost and resource allocation: Data warehousing implementation can be resource-intensive, requiring investments in hardware, software, skilled personnel, and ongoing maintenance. Organizations need to carefully consider the costs involved, including infrastructure costs,

licensing fees, and personnel expenses, to ensure that the benefits of data warehousing outweigh the investments made.

❖ Addressing these implementation issues requires careful planning, collaboration between business and IT teams, and a focus on continuous improvement. Organizations should thoroughly assess their requirements, develop a comprehensive implementation strategy, and regularly monitor and refine their data warehousing processes to maximize the value derived from their data assets.

### 3.4.1 Hardware and software requirements for data warehousing

Hardware and software requirements for data warehousing can vary depending on the specific needs and scale of an organization's data warehousing initiative.

*Hardware Requirements:*

1.  Server Infrastructure: A robust server infrastructure is necessary to host the data warehouse. This includes powerful servers with sufficient processing power, memory (RAM), and storage capacity to handle the data volume and query workloads. The infrastructure may consist of multiple servers, such as database servers, application servers, and storage servers, depending on the architecture and scalability requirements.

2.  Storage Systems: Data warehousing involves storing large volumes of data, so organizations need to ensure sufficient storage capacity. High-performance storage systems, such as RAID arrays or solid-state drives (SSDs), are commonly used to provide fast access to data. Additionally, organizations may consider options like network-attached storage (NAS) or storage area network (SAN) for centralized and scalable storage solutions.

3.  Backup and Recovery: Implementing a robust backup and recovery strategy is critical to protect the data warehouse from data loss or system failures. This may involve regular backups to external storage, replication across multiple servers, or implementing backup and recovery tools provided by the database management system (DBMS) being used.

*Software Requirements:*

1.  Database Management System (DBMS): A reliable and scalable DBMS is a core component of a data warehousing solution. Popular choices for data warehousing include Oracle Database, Microsoft SQL Server, IBM Db2, and PostgreSQL. The DBMS should support features like parallel processing, partitioning, indexing, and query optimization to handle large datasets and complex queries efficiently.

2. Extract, Transform, Load (ETL) Tools: ETL tools are essential for data integration and transformation processes in data warehousing. These tools facilitate extracting data from various source systems, performing data cleansing and transformation, and loading it into the data warehouse. Examples of popular ETL tools include Informatica PowerCenter, IBM InfoSphere DataStage, Microsoft SQL Server Integration Services (SSIS), and Apache NiFi.

3. Business Intelligence (BI) Tools: BI tools enable users to visualize and analyze data stored in the data warehouse. These tools provide features like reporting, dashboards, data visualization, and ad-hoc querying. Popular BI tools include Tableau, Power BI, MicroStrategy, and Looker.

4. Data Modeling and Design Tools: Data modeling tools help in designing the structure and relationships of the data warehouse schema. They provide functionalities to create conceptual, logical, and physical data models, as well as perform schema management. Examples include ER/Studio, Toad Data Modeler, and Oracle SQL Developer Data Modeler.

5. Security and Access Control Tools: Data warehousing requires robust security measures to protect sensitive data. Access control tools and mechanisms, such as user authentication, role-based access control (RBAC), and encryption, should be implemented. Database-specific security features and third-party security tools can assist in securing the data warehouse.

---

*Check your progress 2*

*1. What are the various hardware requirements for implementing a Data Warehouse.*

*2. What are the major issues in the implementation of data warehouse?*

*3. What is backup and recovery important?*

*4. Which type of storage devices are commonly used to provide fast access?*

*5. What kind of features DBMS should support to handle large datasets?*

---

## 3.6 Summary

In this unit we focused on the essential components of a Data warehouse. The 3-layered architecture of Data warehouse and various implementation issues related to Data warehouse were explained in this unit. The unit explained the following:

- ➢ Data warehousing components
- ➢ 3-level architecture of Data warehouse
- ➢ Implementation considerations of Data warehouse
- ➢ Hardware and Software requirements for implementing Data warehouse.

## 3.7 Review Questions

Q1. Discuss 3-level architecture of Data Warehouse. Give a diagram for the 3-level architecture.

_____

_____

_____

Q2. Give a diagram for the essential modules of a Data Warehouse. Explain the functionality of each of the modules.

_____

_____

_____

Q3. Discuss the essential implementation issues of Data Warehouse.

_____

_____

_____

Q4. Explain the essential hardware and software requirement to implement a Data warehouse.

_____

_____

_____

Q5. Explain the various implementation issues related to a data warehouse.

_____

_____

_____

# Suggested Further Readings

1. "Data Mining: Concepts and Techniques" by Jiawei Han, Micheline Kamber, and Jian Pei: This book provides a comprehensive introduction to the concepts, methods, and algorithms of data mining. It covers various data mining techniques, including classification, clustering, association analysis, and outlier detection.

2. "Data Mining: Practical Machine Learning Tools and Techniques" by Ian H. Witten, Eibe Frank, and Mark A. Hall: This book offers practical guidance on using data mining techniques in real-world applications. It covers a wide range of data mining methods, algorithms, and tools, with a focus on applying them using the popular open-source software, Weka.

3. "Introduction to Data Mining" by Pang-Ning Tan, Michael Steinbach, and Vipin Kumar: This book provides a comprehensive introduction to the fundamental concepts and techniques of data mining. It covers topics such as data preprocessing, classification, clustering, association analysis, and anomaly detection.

4. "Data Mining: Concepts, Models, Methods, and Algorithms" by Mehmed Kantardzic: This book presents a comprehensive overview of data mining concepts, models, methods, and algorithms. It covers key topics such as data preprocessing, classification, clustering, pattern mining, and data visualization.

**Master of Computer Application**
# MCA-E6N
# Data Mining

उ० प्र० राजर्षि टण्डन
मुक्त विश्वविद्यालय, प्रयागराज

# Block
# 2

## Data Mining and its Techniques

## Course Design Committee

| | |
|---|---|
| **Prof. Ashutosh Gupta** | **Chairman** |
| Director (In-charge) | |
| School of Computer & Information Science, UPRTOU Allahabad | |
| **Prof. Suneeta Agarwal** | **Member** |
| Dept. of Computer Science & Engineering | |
| Motilal Nehru National Institute of Technology Allahabad | |
| **Dr. Upendra Nath Tripathi** | **Member** |
| Associate Professor | |
| DeenDayalUpadhyay Gorakhpur University, Gorakhpur | |
| **Dr. Ashish Khare** | **Member** |
| Associate Professor | |
| Dept. of Computer Science, University of Allahabad, Prayagraj | |
| **Ms. Marisha** | **Member** |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |
| **Mr. Manoj Kumar Balwant** | **Member** |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |

## Course Preparation Committee

| | |
|---|---|
| **Dr. Tulika Narang** | **(Block 1, 2 & 4) Author** |
| Assistant Professor, Computer science | |
| United University, Rawatpur, Prayagraj | |
| **Dr. Krishan Kumar** | **(Block 3-Unit 7) Author** |
| Assistant Professor, | |
| Department of Computer Science, Faculty of Technology | |
| Gurukula Kangri Vishwavidyalaya, Haridwar (UK) | |
| **Dr. Pooja Yadav** | **(Block 3-Unit 8) Author** |
| Assistant, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Iram Naim** | **(Block 3-Unit 9) Author** |
| Assistant Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Brajesh Kumar** | **Editor** |
| Associate Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Mr. Manoj Kumar Balwant** | **Coordinator** |
| Assistant Professor (Computer Science), | |
| School of Sciences, UPRTOU, Prayagraj | |

# Block-2 Introduction

In this block, the fourth unit explores the concept of data mining, emphasizing its significance in extracting valuable insights from large datasets. *Data mining* is a multifaceted process that utilizes various techniques and algorithms to uncover patterns, relationships, and trends that inform decision-making and knowledge discovery. Its applications span various fields, including business, marketing, finance, healthcare, telecommunications, e-commerce, and scientific research, enabling organizations to make data-driven decisions and uncover hidden insights.

In the fifth unit, we will explore the specific tasks involved in the data mining. Essential tasks such as association mining, clustering, classification, and deviation detection will be examined in detail. This unit will also cover various clustering methods and provide a comparative analysis of classification versus clustering, along with different approaches to deviation detection.

The sixth unit introduces advanced concepts such as neural networks, genetic algorithms, and rough sets. We will discuss the structure of neural networks, their data mining applications, and the back propagation algorithm's **complexity**. Inspired by the human brain's architecture, Neural networks are particularly adept at handling complex patterns and high-dimensional data, making them invaluable for classification, regression, clustering, and anomaly detection tasks. Additionally, we will explore the principles of genetic algorithms and their functionalities, highlighting their role in optimizing data mining processes.

# UNIT 4: Introduction to Data Mining

**Structure**

## 4.0 Introduction

In this unit, the focus is to understand the concept of Data Mining. The focus of this unit is on various issues of data mining. The essential tasks of data mining are also discussed here. Various metrics usedin data mining is also part of this unit. Data Mining refers to the process of extracting valuable or meaningful information from large datasets. It involves the use of various techniques and algorithms to discover patterns, relationships, and insights that can be used for decision-making, prediction, and knowledge discovery.Data mining is widely used in various fields, including business and marketing, finance, healthcare, telecommunications, e-commerce, and scientific research. It helps organizations make informed decisions, discover hidden patterns, and gain valuable insights from their data.

## 4.1 Objectives
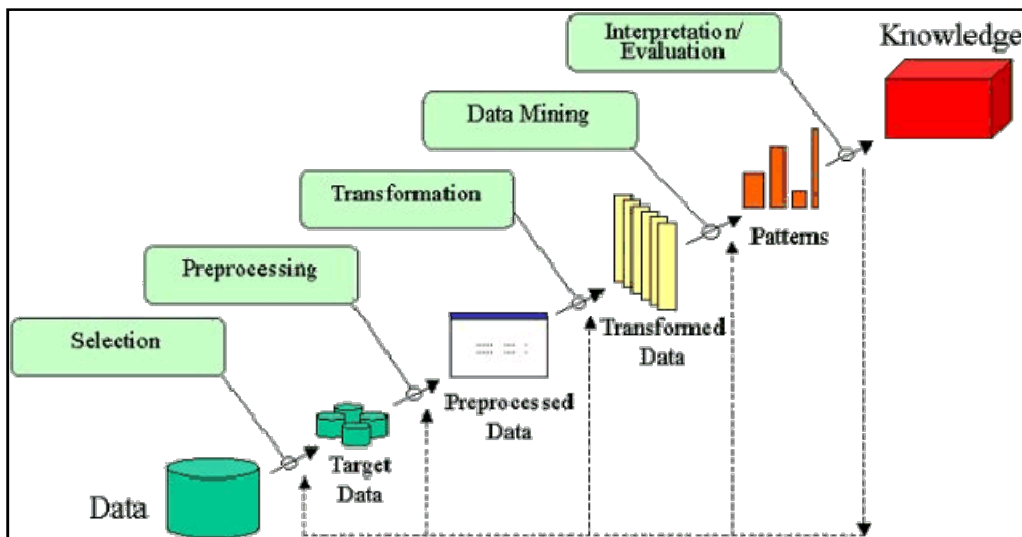
After the end of this unit, you should be able to:

- Explain different tasks of Data Mining process

- Understand Data mining from Database perspective

- Understand and able to define Data Mining metrics

- Identify essential data mining issues

## 4.2 Data Mining Tasks

**Data mining** is the method of finding useful patterns iof knowledge in data. Data mining includes various tasks and methods to extract and discover useful and valuable patterns of knowledge from huge data sets. It the analysis step in Knowledge discovery (KDD) in Databases process. Sometimes Data Mining and Knowledge Discovery in Databases is used as similar keywords. But both are different. Data Mining is a step in Knowledge Discovery in Databases (KDD) process. The steps of Knowledge Discovery in Databases are as follows:

**1. Data cleaning** -to eliminate noise and irrelevant data.

**2. Data integration** - multiple data sources are combined to single repository.

**3. Data selection** – Relevant data from the database is extracted for analysis.

**4. Data transformation** - Data is changed/modified and consolidated into forms suitable for mining. This is done by applying aggregation functions.

**5. Data mining** –In this step an essential process where data mining methods are applied to find hidden valuable patterns of knowledge.

**6. Pattern evaluation** – In this identify the interesting patterns are represented and evaluated using various measures and metrics.

**7. Knowledge presentation-** Data visualization techniques are used to present mined knowledge to users.



Source: https://www2.cs.uregina.ca/~dbd/cs831/notes/kdd/1_kdd.html

Figure 1: Knowledge Discovery in Databases

The essential data mining tasks are classified as:

- ➢ Classification

- ➢ Clustering
- ➢ Association Rule Mining
- ➢ Regression
- ➢ Outlier Analysis

Data mining involves a range of tasks and techniques aimed at discovering patterns, relationships, and insights from large datasets. Here are some common data mining tasks:

1. **Association Rule Mining**: This task involves discovering relationships or associations among items in a dataset. It identifies frequent item sets and generates rules that describe the co-occurrence of items in transactions. For example, in a retail setting, association rule mining can reveal that customers who buy diapers are likely to also purchase baby formula.

2. **Classification**: Classification is the process of assigning predefined categories or labels to instances based on their features or attributes. It involves building predictive models that can learn from labeled training data and then classify new, unseen instances. For instance, classifying emails as spam or non-spam based on their content is a common classification task.

3. **Clustering**: Clustering aims to group similar instances together based on their characteristics without prior knowledge of the class labels. It helps identify natural groupings or patterns in the data. For example, clustering can be used to segment customers into distinct groups based on their purchasing behavior.

4. **Regression Analysis**: Regression analysis is used to model and analyze the relationship between a dependent variable and one or more independent variables. It helps predict continuous numeric values. For instance, predicting house prices based on features like location, size, and number of rooms is a regression task.

5. **Anomaly Detection**: Anomaly detection focuses on identifying unusual or abnormal instances in a dataset. It aims to uncover patterns or behaviors that deviate significantly from the norm. Anomaly detection is valuable in various domains, such as fraud detection, network intrusion detection, and equipment failure prediction.

6. **Text Mining**:

   Text mining, also known as text analytics or natural language processing (NLP), is a data mining technique that focuses on extracting meaningful information and insights from unstructured textual data. It involves applying computational methods to process, analyze, and interpret text in order to discover patterns, sentiments, relationships, and other valuable knowledge.Text mining involves extracting meaningful information and insights from unstructured text data. It

encompasses tasks like text categorization, sentiment analysis, topic modeling, and entity recognition. Text mining enables the analysis of large volumes of text from sources like social media, customer reviews, and documents.

7. **Time Series Analysis**: Time series analysis deals with data collected over time, typically at regular intervals. It focuses on understanding patterns, trends, and dependencies in temporal data. Time series analysis is useful for forecasting future values, detecting anomalies, and analyzing sequential patterns.

8. **Dimensionality Reduction**: Dimensionality reduction techniques aim to reduce the number of variables or features in a dataset while preserving important information. They help overcome challenges associated with high-dimensional data and can improve model efficiency and interpretability.

9. **Outlier Analysis**: Outlier nalysis, also known as outlier detection or anomaly detection, is a data mining technique that focuses on identifying observations or data points that deviate significantly from the normal behavior or patterns of the dataset. Outliers are data points that are considerably different from the majority of the data and may indicate errors, anomalies, or interesting events in the dataset. Outlier analysis is applied in various domains such as fraud detection, network intrusion detection, sensor data analysis, medical diagnostics, quality control, and financial market analysis. By identifying unusual patterns or behaviors, outlier analysis helps in uncovering anomalies, detecting potential risks or threats, and improving decision-making processes.

---

*Check your progress 1*

*1. What is data mining? Discuss atleast 4 essential tasks of Data Mining.*

*2. Differentiate between classification and clustering.*

*3. What is knowledge discovery in databases?*

*4. List major steps in KDD?*

*5. List at least 4 essential tasks of Data Mining.*

---

# 4.3 DATA MINING ISSUES

Data mining, while a powerful tool for extracting valuable insights from large datasets, is not without its challenges and issues. Some of the common issues in data mining include:

1. **Data Quality:** Data mining relies heavily on the quality and reliability of the underlying data. If the data is incomplete, inconsistent, or contains errors, it can significantly affect the accuracy and validity of the mining results. Data cleaning and preprocessing techniques are often required to address these issues.

2. **Data Privacy and Security:** Data mining involves accessing and analyzing large amounts of data, which can raise privacy concerns. Sensitive information, such as personally identifiable information (PII), needs to be handled carefully to ensure compliance with privacy regulations and to protect individuals' privacy rights.

3. **Data Integration:** In many cases, data used for mining is collected from various sources and may be stored in different formats and structures. Integrating and combining these diverse datasets can be challenging and time-consuming, as it requires resolving inconsistencies and dealing with varying levels of data granularity.

4. **Dimensionality and Complexity:** Data mining deals with datasets that can have a high number of variables or features, known as high dimensionality. High-dimensional data poses challenges in terms of computational efficiency and interpretability of results. Additionally, complex relationships and interactions among variables can make it difficult to extract meaningful patterns and insights.

5. **Overfitting and Model Selection:** Over fitting occurs when an algorithm fits too closely or even exactly to its training data, resulting in a model that can't make accurate predictions or conclusions from any data other than the training data. A statistical model is said to be overfitted when the model does not make accurate predictions on testing data. When a model gets trained with so much data, it starts learning from the noise and inaccurate data entries in our data set. And when testing with test data results in High variance. Then the model does not categorize the data correctly, because of too many details and noise. The causes of overfitting are the non-parametric and non-linear methods because these types of machine learning algorithms have more freedom in building the model based on the dataset and therefore they can really build unrealistic models. A solution to avoid overfitting is using a linear algorithm if we have linear data or using the parameters like the maximal depth if we are using decision trees. When building predictive models through data mining, there is a risk of overfitting, where the model

performs well on the training data but fails to generalize to new, unseen data. Selecting the appropriate modeling techniques and avoiding overfitting require careful consideration and validation.

6. **Interpretability and Explainability:** Some data mining techniques, such as deep learning and complex ensemble models, can be highly accurate but lack interpretability. Understanding and explaining the reasoning behind the mining results can be challenging, particularly when dealing with complex algorithms.

7. **Ethical Considerations:** Data mining can raise ethical concerns, especially when it involves sensitive or personal data. Ensuring fairness, transparency, and accountability in data mining practices is crucial to avoid biases and discriminatory outcomes.

8. **Scalability:** As the size of datasets continues to grow exponentially, scalability becomes a significant issue in data mining. Efficient algorithms and computational resources are necessary to handle large-scale data mining tasks within reasonable time frames.

## 4.4 Data Mining Metrics

Data mining metrics are quantitative measures used to evaluate the performance and effectiveness of data mining models and algorithms. These metrics provide insights into the accuracy, quality, and reliability of the mining results. The metrics consider a confusion matrix structure for evaluation. A confusion matrix provides a tabular representation of the predicted and actual class labels. It shows the number of true positives, true negatives, false positives, and false negatives. From the confusion matrix, various metrics like accuracy, precision, recall, and F1 score can be derived. A confusion matrix is a table used to evaluate the performance of a classification model. It provides a comprehensive view of how well the model is predicting the actual classes. The matrix consists of four key components:

1. **True Positives (TP)**: The number of instances correctly predicted as positive.
2. **True Negatives (TN)**: The number of instances correctly predicted as negative.
3. **False Positives (FP)**: The number of instances incorrectly predicted as positive (Type I error).
4. **False Negatives (FN)**: The number of instances incorrectly predicted as negative (Type II error).

For a binary classification problem, the confusion matrix is typically structured as follows:

|  | Predicted Positive | Predicted Negative |
|---|---|---|
| Actual Positive | True Positive (TP) | False Negative (FN) |
| Actual Negative | False Positive (FP) | True Negative (TN) |

Example

Suppose we have a binary classification model that predicts whether an email is spam or not spam. After testing the model, we get the following results:

- 50 emails are correctly identified as spam (TP)
- 40 emails are correctly identified as not spam (TN)
- 10 emails are incorrectly identified as spam (FP)
- 5 emails are incorrectly identified as not spam (FN)

The confusion matrix would look like this:

|  | Predicted Spam | Predicted Not Spam |
|---|---|---|
| Actual Spam | 50 | 5 |
| Actual Not Spam | 10 | 40 |

Some commonly used data mining metrics are:

**1.Accuracy**: Accuracy measures the overall correctness of a predictive or classification model. It is typically calculated as the ratio of correctly predicted instances to the total number of instances. High accuracy indicates that the model is making accurate predictions.Accuracy is a key metric and is used to evaluate the performance of a classification model. It measures the proportion of correctly classified instances out of the total instances. In other words, it gives the overall effectiveness of the model in correctly identifying both positive and negative instances.The formula for accuracy is:

$$\frac{TP + TN}{TP + TN + FP + FN}$$

where:

- *TP*: True Positives (instances correctly predicted as positive)
- *TN*: True Negatives (instances correctly predicted as negative)
- *FP*: False Positives (instances incorrectly predicted as positive)
- *FN*: False Negatives (instances incorrectly predicted as negative)

Example-

Consider a binary classification problem where the confusion matrix is as follows:

|  | **Predicted Positive** | **Predicted Negative** |
|---|---|---|
| **Actual Positive** | 50 | 5 |
| **Actual Negative** | 10 | 40 |

Using the formula for accuracy:

Accuracy=$\frac{50+40}{50+40+10+5}=\frac{90}{105}$=0.857. This implies that the accuracy of the model is

approximately 85.7%.

**2. Precision and Recall**: Precision and recall are metrics used in binary classification tasks. Precision measures the proportion of true positive instances among the predicted positive instances. Recall, also known as sensitivity, measures the proportion of true positive instances correctly identified. A balance between precision and recall is important, as optimizing one may affect the other.PrecisionandRecall are two important metrics used to evaluate the performance of a classification model, particularly in situations where class imbalance is present. They help in understanding how well the model is predicting the positive class and how effective it is in identifying all relevant instances of the positive class.Precision is also known as Positive Predictive Value. It measures the proportion of true positive predictions among all positive predictions made by the model.

The precision is determined by the following equation:
$$Precision=\frac{TP}{TP+FP}$$

where:

- *TP* **(True Positives)**: Instances correctly predicted as positive.
- *FP* **(False Positives)**: Instances incorrectly predicted as positive.

The parameter recall is calculated as follows:

$$Recall=\frac{TP}{TP+FN}$$

where:

- **TP (True Positives)**: Instances correctly predicted as positive.
- **FN (False Negatives)**: Instances incorrectly predicted as negative.

Let us illustrate the calculations of Precision and Recall with the help of an example. Consider the following confusion matrix.Consider the confusion matrix:

|  | Predicted Positive | Predicted Negative |
|---|---|---|
| Actual Positive | 50 | 5 |
| Actual Negative | 10 | 40 |

By applying the formula, the Precision is obtained as follows.

Precision=$\frac{TP}{TP+FP}$=$\frac{50}{50+10}$=0.833. Thus the precision of the model is approximately 83.3%.

Using the formula for recall:

Recall= $\frac{TP}{TP+FN}$=$\frac{50}{50+5}$ =0.909.

Thus the recall of the model is approximately 90.9%

**3. F1 Score**: The F1 score is a combined metric that considers both precision and recall. It is the harmonic mean of precision and recall and provides a single measure to evaluate the model's performance. F1 score is useful when there is an imbalance between the positive and negative classes. The F1 Score is the harmonic mean of precision and recall, providing a single metric that balances both concerns. It is particularly useful when you need a single metric to compare models or when dealing with imbalanced classes.

Formula for F1 Score

$$2 * \frac{Precision * Recall}{Precision + Recall}$$

=2 *$\frac{0.833*0.909}{0.833+0.909}$= 0.87.  Thus the F1 Score of the model is approximately 87%.

## 4.5 Data Mining from a Database perspective

From a database perspective, data mining involves the extraction of useful patterns, trends, and knowledge from large datasets stored in databases. Some key considerations and techniques related to data mining from a database perspective are:

1) **Data Preparation**: Before data mining can take place, the relevant data needs to be extracted, transformed, and loaded (ETL) from the database. This involves selecting the appropriate tables or views, handling missing or inconsistent data, and transforming the data into a suitable format for analysis.

2) **Data Integration**: Data mining often requires combining data from multiple databases or data sources. Data integration involves merging and reconciling data from different sources into a unified view. This process may involve resolving schema and data format differences, dealing with data duplication, and ensuring data consistency.

3) **Database Querying**: To perform data mining tasks, queries are executed against the database to retrieve the relevant data. SQL (Structured Query Language) is commonly used for querying relational databases. Complex queries may involve joins, aggregations, and filtering to extract the necessary data subsets for analysis.

4) **Indexing and Performance Optimization**: Efficient indexing and query optimization techniques are crucial for data mining tasks. Indexes help speed up query execution by providing faster access to specific data subsets. Techniques like query optimization, caching, and parallel processing can significantly improve the performance of data mining operations.

5) **Sampling and Partitioning**: When dealing with large databases, it may not be feasible or necessary to analyze the entire dataset. Sampling techniques can be employed to extract

representative subsets of the data for analysis. Partitioning techniques, such as horizontal or vertical partitioning, can divide the data into smaller, manageable chunks for efficient processing.

6) **Data Security and Privacy**: Databases often contain sensitive or confidential information, and data mining should adhere to security and privacy regulations. Access control mechanisms, encryption, and anonymization techniques may be implemented to protect the data during the mining process and ensure compliance with privacy laws.

7) **Incremental and Stream Mining**: In some cases, data mining needs to be performed on continuously incoming data or in real-time. Incremental and stream mining techniques enable the analysis of data as it arrives, allowing for timely insights and decision-making.

8) **Data Warehousing**: Data mining can leverage data warehouses, which are specialized databases designed for supporting analytics and decision-making. Data warehouses provide a consolidated view of data from various sources, optimized for querying and analysis.

Database technologies and techniques play a crucial role in supporting efficient and effective data mining. They provide the foundation for data extraction, integration, querying, and storage, enabling the application of various data mining algorithms and techniques to uncover valuable insights from large-scale databases.

## 4.6 Essential Aspects In Data Mining

Various essential aspects are related to the process of data mining.Various aspects such as defining the problem, data preparation, model selection, training and evaluation are considered for a data mining problem. The key aspects of data mining include:

1. **Problem Definition:** Clearly defining the problem or objective to be addressed through data mining. This involves understanding the business goals, identifying the relevant data sources, and determining the specific questions to be answered or insights to be gained. Problem definition is the first step and an essential step in defining the problem and identifying the goals.

2. **Data Preparation:** Collecting, cleaning, integrating, and transforming the raw data into a suitable format for analysis. This includes handling missing values, resolving inconsistencies, removing duplicates, and performing data normalization or scaling.
   a. Various methods of data data preparation such as binning by means, clustering and outlier detection are used in this step.

3. **3) Exploratory Data Analysis (EDA):** Exploring and understanding the data through statistical techniques and visualizations. EDA helps identify patterns, trends, distributions, and relationships

within the data, which can guide the subsequent data mining steps.Exploratory Data Analysis (EDA) is a crucial step in the data analysis process, where you use statistical and graphical techniques to understand the structure, patterns, and relationships in a dataset. The goal of EDA is to gain insights, detect anomalies, test hypotheses, and check assumptions before applying more sophisticated modeling techniques.

4. **Feature Selection and Engineering:** Selecting the most relevant features (variables) from the dataset that are likely to have a significant impact on the outcome. Feature engineering involves creating new features or transforming existing ones to improve the performance of the data mining models.

5. **Model Selection:** Choosing appropriate data mining models or algorithms that are well-suited for the specific problem and data characteristics. This depends on factors such as the type of data (e.g., numerical, categorical, text), the goal of analysis (e.g., prediction, classification, clustering), and available computational resources.

6. **Model Training and Evaluation:** Training the selected models using labelled or historical data and evaluating their performance. This includes techniques such as cross-validation, training/testing splits, and performance metrics (e.g., accuracy, precision, recall, F1 score) to assess the quality and predictive power of the models.

7. **Model Interpretation and Validation:** Interpreting the results of the data mining models to gain insights and validate their usefulness. This involves understanding the underlying patterns or rules discovered by the models and assessing their validity and reliability.

8. **Deployment and Implementation:** Translating the insights and findings from the data mining process into actionable recommendations or solutions. This may involve integrating the models into existing systems or processes, developing dashboards or reports for stakeholders, or implementing automated decision-making systems based on the mining results.

9. **Monitoring and Maintenance:** Continuously monitoring the performance of the deployed models and updating them as new data becomes available or the problem requirements change. This ensures that the models remain accurate and effective over time.

## 4.7 Summary

In this unit we presented various tasks in Data Mining process. We also described the data mining process from database perspective. Various aspects used in data mining and various data mining tasks were also discussed in this block. To summarize the unit covers the following:

- ✓ Basic Data Mining Tasks
- ✓ Data Mining Issues
- ✓ Data Mining Metrics
- ✓ Data Mining from a Database Perspective
- ✓ Data Mining essential aspects

## REVIEW QUESTIONS

1. Discuss Data integration and Data Warehouse.

_____

_____

2. Discuss Knowledge Discovery in Databases.

_____

_____

3. Explain Precision and Recall metrics.

_____

_____

4. What are outliers and what is outlier analysis?

_____

_____

5. Explain essential data mining tasks.

_____

# UNIT 5: Data Mining Techniques

**Structure**

## 5.0 Introduction

In previous unit of this block, you learnt various tasks in Data Mining process. The data mining process is a structured approach used to extract meaningful patterns, relationships, and insights from large datasets. It typically involves several key stages, each with specific tasks designed to prepare data, apply algorithms, and interpret results. The data mining process from database perspective was also explained in the previous unit.The various aspects used in data mining and various data mining tasks were also discussed in the previous block. This block focuses and covers the following various essential tasks of Data Mining such as Association rule mining, clustering, classification and Deviation detection.

## 5.1 Objectives

This unit deals with the representation of algorithm in terms of flowchart and pseudo code.

At the end of this unit, you will be able to:

- Understand and Explain Association Rule Mining

- Understand Clustering and various types

- Understand Classification and comparison with clustering

- Understand and Explain Deviation Detection

## 5.2 Association rule mining

Association rule mining is a technique used in data mining and machine learning to discover interesting relationships, patterns, and associations among items in large datasets. It aims to find associations between items based on their co-occurrence in transactions or events.

The most commonly used algorithm for association rule mining is the Apriori algorithm. The Apriori algorithm works by scanning the dataset multiple times and generating candidate item sets of increasing length. It then calculates the support and confidence measures for each itemset to determine the significant associations.The key terms and concepts associated with association rule mining:

1) Itemset: An itemset is a collection of items that appear together in a transaction or event. It can be a single item (singleton) or a combination of multiple items.

2) Support: Support measures the frequency or occurrence of an itemset in the dataset. It indicates the proportion of transactions in which an itemset appears.

3) Confidence: Confidence measures the strength of an association rule. It represents the conditional probability of finding the consequent item(s) in a transaction given the presence of the antecedent item(s).

4) Association Rule: An association rule is a statement that captures the relationship between items. It consists of an antecedent (premise) and a consequent (outcome) separated by an arrow. For example, "if {A} then {B}" denotes that there is a relationship between items A and B.

5) Support Threshold: It is a user-defined parameter that specifies the minimum support required for an itemset or association rule to be considered significant. It helps filter out less frequent or less interesting associations.

6) Lift: Lift is a measure of the strength of an association rule. It compares the observed support of the rule to the expected support if the antecedent and consequent were independent of each other. A lift greater than 1 indicates a positive association, while a lift less than 1 suggests a negative association.

Association rule mining has numerous applications, including market basket analysis, recommendation systems, web usage mining, customer behavior analysis, and more. It helps identify frequently co-occurring items, understand buying patterns, and make data-driven decisions based on discovered associations.Association rule mining is a powerful technique but computationally expensive for large datasets. Several optimized algorithms and techniques have been developed to improve efficiency and scalability, such as FP-growth.

## 5.2.1 APRIORI ALGORITHM

The Apriori algorithm is a classic algorithm used for association rule mining, specifically for finding frequent itemsets in a dataset. It was proposed by Rakesh Agrawal and Ramakrishnan Srikant in 1994. *The Apriori algorithm is based on the principle of "apriori property," which states that if an itemset is frequent, then all of its subsets must also be frequent.*

The Apriori algorithm works as follows:

1) Generating frequent 1-itemsets: The algorithm starts by scanning the dataset to determine the frequency of each item (single itemset). It counts the occurrences of each item and identifies the frequent items based on a user-defined minimum support threshold.

2) Generating candidate $k$-itemsets: The algorithm generates candidate itemsets of length $k$ ($k > 1$) by combining frequent ($k$-1)-itemsets. It employs the apriori property to ensure that only promising itemsets are considered. The candidates are formed by taking the union of frequent (k-1)-itemsets with the same prefix.

3) Scanning the dataset: The algorithm scans the dataset again to determine the frequency of each candidate itemset. It counts the occurrences of each candidate in the transactions and identifies the frequent itemsets based on the minimum support threshold.

4) Iterative process: Steps 2 and 3 are repeated iteratively until no more frequent itemsets can be found. In each iteration, the algorithm generates candidate (k+1)-itemsets based on the frequent k-itemsets discovered in the previous step and scans the dataset to identify the frequent itemsets.

5) Generating association rules: Once the frequent itemsets are discovered, the algorithm generates association rules from the frequent itemsets. It uses a user-defined minimum confidence threshold to filter out uninteresting or weak associations. Association rules are generated by considering all possible combinations of items within the frequent itemsets.

The Apriori algorithm can be computationally expensive, especially for large datasets, because it requires multiple passes over the data and generates a large number of candidate itemsets. To improve efficiency, optimization techniques such as pruning, hashing, and tree-based data structures like the FP-tree (used in the FP-growth algorithm) can be employed.The Apriori algorithm is widely used for market basket analysis, where it helps identify frequently co-occurring items in customer transactions. It has also been applied to various other domains, such as web mining, bioinformatics, and text mining.

**Example:**

Let's consider a retail scenario where association rule mining can be applied. Suppose we have a transaction dataset that records customer purchases in a grocery store. Each transaction consists of a set of items that a customer bought during a single visit.

Here's a simplified example of the transaction dataset:

Transaction 1: {Bread, Milk, Eggs}

Transaction 2: {Bread, Diapers}

Transaction 3: {Milk, Diapers, Beer}

Transaction 4: {Bread, Milk, Diapers, Beer}

Transaction 5: {Bread, Milk, Diapers}

Apply association rule mining to discover interesting relationships between items. One commonly used metric in association rule mining is called support, which measures the frequency of a particular itemset or combination of items in the dataset.

Let's say we want to find association rules with a minimum support of 40% and a minimum confidence of 60%. This means we are interested in finding itemsets that appear together in at least 40% of the transactions and rules that have a confidence level of at least 60%.

Based on the dataset, we can calculate the support and confidence for different itemsets and rules. Here are a couple of example rules that could be discovered:

Rule 1: {Milk} => {Bread} Support: 60% (appears in 3 out of 5 transactions) Confidence: 75% (3 transactions with Milk and 4 transactions with Bread)

This rule indicates that if customer purchases milk, they are likely to buy bread as well. The support of 60% suggests that milk and bread appear together in 60% of the transactions, while the confidence of 75% indicates that out of all transactions with milk, 75% of them also contain bread.

Rule 2: {Diapers} => {Beer} Support: 40% (appears in 2 out of 5 transactions) Confidence: 66.7% (2 transactions with Diapers and 3 transactions with Beer)

This rule suggests that customers who buy diapers also tend to purchase beer. The support of 40% indicates that diapers and beer appear together in 40% of the transactions, while the confidence of 66.7% suggests that out of all transactions with diapers, 66.7% of them also include beer.

These association rules provide insights into the buying patterns of customers in the grocery store and can be used for various purposes, such as product placement, targeted promotions, or cross-selling strategies.

By leveraging association rule mining, retailers can optimize their sales and marketing strategies based on the discovered associations between items.

## 5.3 Clustering

Clustering is a machine learning technique used to group similar objects or data points into clusters based on their inherent similarities or patterns. It is an unsupervised learning method because it does not rely on predefined labels or target variables. Clustering algorithms aim to discover underlying structures or relationships within the data by grouping similar items together and separating dissimilar items.The main goal of clustering is to maximize the intra-cluster similarity (similarity among objects within a cluster) and minimize the inter-cluster similarity (similarity between objects in different clusters). Clustering can be applied to various types of data, such as numerical data, categorical data, text data, and more.Some commonly used clustering algorithms are:

1) K-means: K-means is a popular partition-based clustering algorithm. It starts by randomly selecting K initial cluster centers and assigns each data point to the nearest cluster center based on a distance metric (usually Euclidean distance). The cluster centers are then updated iteratively by computing the mean of the data points in each cluster. This process continues until convergence, where the cluster assignments and cluster centers stabilize.The steps of the algorithm are:

1. **Initialization**:
   - Select k, the number of clusters.
   - Randomly initialize k cluster centroids. These can be randomly selected data points or chosen using more sophisticated methods like K-means++.
2. **Assignment Step**:
   - Assign each data point to the nearest cluster centroid. This is typically done using the Euclidean distance formula, but other distance metrics can be used.
3. **Update Step**:
   - Recalculate the centroids of the clusters by taking the mean of all data points assigned to each cluster.
4. **Repeat**:
   - Repeat the assignment and update steps until the centroids no longer change significantly or a maximum number of iterations is reached.

The detailed steps of K-means with exampleis as follows:

1. **Initialization**:
   - Suppose we have a dataset with points in a 2D space and we want to form 3 clusters (k=3).
   - Randomly select 3 initial centroids, C1, C2, and C3.
2. **Assignment Step**:
   - For each data point, calculate the distance to each centroid and assign the point to the closest centroid.
   - Example: If point P1 is closest to centroid C2, assign P1 to cluster 2.
3. **Update Step**:
   - Recalculate the centroids of each cluster by finding the mean position of all points assigned to that cluster.
   - Example: If cluster 2 contains points P1, P4, and P7, the new centroid C2 is the mean of these points.
4. **Repeat**:
   - Continue the assignment and update steps until the centroids stabilize.

**The advantages of the algorithm are:**

1. **Simplicity**: Easy to understand and implement.
2. **Efficiency**: Computationally efficient with a time complexity of O(n·k.d) where n is the number of data points, k is the number of clusters, iii is the number of iterations, and d is the number of dimensions.
3. **Scalability**: Works well with large datasets

**The limitations of the algorithm are**:

1. **Choice of k**: Requires the number of clusters k to be specified in advance.
2. **Sensitivity to Initialization**: The final clusters can vary based on the initial choice of centroids.
3. **Assumption of Spherical Clusters**: Assumes clusters are spherical and of similar size, which may not always be the case.
4. **Sensitivity to Outliers**: Outliers can significantly affect the cluster centroids.

**Some applications of the algorithm are:**

1. **Market Segmentation**: Grouping customers based on purchasing behavior.

2. **Image Compression**: Reducing the number of colors in an image by clustering similar colors.

3. **Document Clustering**: Grouping similar documents for topic identification.

4. **Anomaly Detection**: Identifying outliers in data, such as fraud detection.

2) Hierarchical Clustering: Hierarchical clustering builds a hierarchy of clusters using a bottom-up (agglomerative) or top-down (divisive) approach. In agglomerative clustering, each data point initially represents a separate cluster, and pairs of clusters are merged iteratively based on a similarity measure until a single cluster containing all data points is formed. In divisive clustering, the process starts with a single cluster containing all data points, and it is recursively divided into smaller clusters until individual data points form their own clusters.

3) Density-Based Spatial Clustering of Applications with Noise (DBSCAN): DBSCAN is a density-based clustering algorithm. It groups together data points that are densely packed, based on a density criterion. It identifies core points (data points with a sufficient number of neighboring points) and expands clusters by including nearby points within a specified radius. It can discover clusters of arbitrary shape and is robust to noise and outliers.

4) Gaussian Mixture Models (GMM): GMM is a probabilistic clustering algorithm that assumes the data points are generated from a mixture of Gaussian distributions. It models the data as a collection of Gaussian components and estimates the parameters of these components to determine the underlying clusters. GMM assigns probabilities to each data point belonging to each cluster, allowing soft assignment of data points to multiple clusters.

5) DBSCAN: DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a density-based clustering algorithm. It groups together data points that are denselypacked, based on a density criterion. It identifies core points (data points with a sufficient number of neighboring points) and expands clusters by including nearby points within a specified radius. It can discover clusters of arbitrary shape and is robust to noise and outliers.

The choice of clustering algorithm depends on the nature of the data, the desired outcome, and the specific requirements of the problem at hand. Each algorithm has its strengths and weaknesses, and different algorithms may perform better on different types of data or clustering scenarios. Clustering has a wide range of applications across various domains, including customer segmentation, image segmentation, anomaly detection, document clustering, and social network analysis, among others.Clustering algorithms can be broadly classified into several categories based on their approach, assumptions, and characteristics. classification of clustering algorithms is:

1.  **Partition-based Clustering:**

    ➢ K-means: Divides the data into K clusters by minimizing the sum of squared distances between data points and the cluster centroid.

    ➢ K-medoids: Similar to K-means, but uses medoids (actual data points) as cluster representatives instead of centroids.

    ➢ Fuzzy C-means: Assigns a degree of membership to each data point for all clusters, allowing soft assignment of points to clusters.

2.  **Hierarchical Clustering:**

    ➢ Agglomerative: Starts with each data point as an individual cluster and merges clusters iteratively based on a distance or similarity measure.

    ➢ Divisive: Begins with all data points in a single cluster and recursively splits clusters until each data point forms its own cluster.

3.  **Density-based Clustering:**

    ➢ DBSCAN: Groups data points based on density, identifying dense regions separated by sparser areas. It can handle clusters of arbitrary shape and detect outliers as noise points.

    ➢ OPTICS: Similar to DBSCAN but produces a density-based ordering of the data points, allowing flexibility in determining cluster boundaries and discovering clusters of varying densities.

4.  **Model-based Clustering:**

    ➢ Gaussian Mixture Models (GMM): Assumes that the data points are generated from a mixture of Gaussian distributions. It estimates the parameters of the distributions to identify clusters.

    ➢ Latent Dirichlet Allocation (LDA): Primarily used for topic modeling in text data, it assumes that documents are a mixture of underlying topics and assigns topics to each document and words to each topic.

5.  **Subspace Clustering:**

    ➢ CLIQUE: Identifies clusters in subspaces of high-dimensional data, accounting for the presence of clusters in different subsets of dimensions.

    ➢ PROCLUS: Discovers clusters in subspaces of multi-dimensional data, combining both partitioning and hierarchical clustering techniques.

1. Explain Association Rule mining as an essential data mining task.

2. Discuss Apriori algorithm used for Association Rule Mining.

3. Write the major applications of association rule mining.

4.What is the utility of Apriori algorithm?

5.What are limitations of Apriori algorithm?

# 5.4 Classification

Classification is a supervised learning technique in machine learning that involves categorizing or assigning predefined labels or classes to input data based on their features or attributes. It is one of the fundamental tasks in supervised learning, where the model learns from a labelled training dataset and then predicts the class labels for unseen or test data.The goal of classification is to build a model that can accurately classify new instances into one of the predefined classes. The process typically involves the following steps:

1. Data Collection: Gathering a labeled dataset where each data instance is associated with a known class label.

2. Data Preprocessing: Cleaning, transforming, and normalizing the data to ensure its quality and consistency. This step may involve handling missing values, removing outliers, or feature scaling.

3. Feature Selection/Extraction: Identifying relevant features or extracting informative representations from the data that can contribute to accurate classification.

4. Model Training: Using the labeled training dataset to train a classification model. The model learns the patterns and relationships between the input features and their corresponding class labels.

5. Model Evaluation: Assessing the performance of the trained model on unseen or test data. Common evaluation metrics for classification include accuracy, precision, recall, F1-score, and area under the receiver operating characteristic (ROC) curve.

6. Model Optimization: Fine-tuning the model's parameters, adjusting hyperparameters, or exploring different algorithms to improve its performance.

7. Prediction: Using the trained and optimized model to make predictions on new, unlabeled instances or data.

Some popular classification algorithms include:

- Logistic Regression

- Decision Trees

- Random Forests

- Support Vector Machines (SVM)

- Naive Bayes

- k-Nearest Neighbors (k-NN)

- Neural Networks (e.g., Feedforward, Convolutional, Recurrent)

The choice of algorithm depends on factors such as the nature of the data, the number of features, the size of the dataset, interpretability requirements, and the complexity of the problem at hand.

Classification has numerous applications across different domains, including spam detection, sentiment analysis, fraud detection, image classification, medical diagnosis, and customer churn prediction, among others.

---

*Check your Progress 2*

Q1. What is clustering?

Q2. Discuss various clustering methods.

Q3. How is clustering different from classification?

---

# 5.5 Deviation Detection

Deviation detection, also known as anomaly detection, is a machine learning technique that focuses on identifying rare or unusual patterns or data points in a dataset that deviate significantly from the normal or expected behaviour. Anomalies are data points that do not conform to the general patterns or distribution of the majority of the data.

The goal of deviation detection is to automatically and accurately identify these anomalies, which can represent potential errors, fraud, intrusions, or other significant events that require attention. Deviation detection can be performed in various domains and applications, including cybersecurity, network monitoring, fraud detection, system health monitoring, and outlier detection.Some common approaches and techniques used for deviation detection:

1. Statistical Methods:

- ❖ Z-Score: Calculates the standard deviation from the mean and identifies data points that are a certain number of standard deviations away from the mean.
- ❖ Gaussian distribution: Assumes that the data follows a normal distribution and uses statistical measures such as mean and standard deviation to identify anomalies.

2. Distance-based Methods:

- ❖ k-Nearest Neighbors (k-NN): Compares the distance between a data point and its k nearest neighbors. If a point has significantly larger distances compared to its neighbors, it is considered an anomaly.
- ❖ Local Outlier Factor (LOF): Measures the local density of a data point compared to its neighbors. Anomalies have significantly lower density compared to the majority of the data points.

3. Clustering Methods:

- ❖ Density-Based Spatial Clustering of Applications with Noise (DBSCAN): Identifies anomalies as data points that do not belong to any cluster or have fewer neighbouring points.
- ❖ Expectation-Maximization (EM): Models the data distribution using a mixture of Gaussian distributions. Anomalies are considered data points with low probability under the model.

4. Machine Learning Techniques:

- ❖ Isolation Forest: Constructs an ensemble of decision trees that isolate anomalies in a data partitioning process.
- ❖ One-Class Support Vector Machines (SVM): Trains a binary classifier on normal data points and aims to find a boundary that separates normal instances from anomalies.

5. Deep Learning Methods:

- ❖ Auto-encoders: Unsupervised neural networks that aim to reconstruct input data. Anomalies result in higher reconstruction errors compared to normal data.
- ❖ Generative Adversarial Networks (GANs): Generate synthetic data based on the learned distribution and identify anomalies by measuring the difference between real and synthetic data.

The choice of deviation detection method depends on the characteristics of the data, the type of anomalies expected, the available labelled data for training (if any), and the specific requirements of the application. Deviation detection is an ongoing process as new data is collected and new types of anomalies can

emerge. Regular monitoring, evaluation, and updating of the deviation detection system are crucial to maintain  its effectiveness over time.

## 5.6 SUMMARY

In this unit we focused on the various tasks in Data Mining process.  The various data mining tasks, Association mining, clustering, classification, deviation detection were also discussed in the block. Various methods of clustering and comparison of classification with clustering is also discussed in this unit. Various approaches of deviation detection are also discussed in this unit

The unit explained the following:
- ✓  .Association Mining
- ✓  Clustering
- ✓  Categorization of clustering methods
- ✓  Classification and phases in classification process
- ✓  Deviation detection

## Review Questions

Q1. Explain Association mining and Apriori algorithm.

_____

_____

Q2. Discuss clustering and essential types of clustering.

_____

_____

Q3. Explain classification as a data mining task.

_____

Q4. Discuss K-means algorithm. What are the advantages and limitations of K-means algorithm.

_____

# UNIT 6: Specialized Data Mining Techniques

**Structure**

## 6.0 Introduction

In this unit, we will see the concepts of Neural network, Genetic algorithms and rough sets. The structure of neural network and its application in data mining is discussed in this unit. Back propagation algorithm is also explained in this unit. The concept of genetic algorithms and their various functions are discussed in this unit.Neural networks are a powerful tool in data mining, particularly for tasks that involve complex patterns, non-linear relationships, and high-dimensional data. They are inspired by the structure and function of the human brain and consist of layers of interconnected "neurons" that process data and learn from it. In data mining, neural networks are used for tasks such as classification, regression, clustering, and anomaly detection.

## 6.1 Objectives

At the end of this unit, you will be able to:

1.  Understand neural network approach to data mining.

2.  Comprehend the role of genetic algorithms in data mining.

3.  Understand the concepts of rough sets.

## 6.2 Neural Network

A neural network is a computational model inspired by the structure and functionality of the human brain. It is composed of interconnected nodes called artificial neurons or "nodes" that work together to process and transmit information. Neural networks are designed to learn and make predictions or decisions from

input data without being explicitly programmed.The basic building block of a neural network is an artificial neuron, also known as a perceptron. Each perceptron takes multiple inputs, applies weights to those inputs, and combines them using an activation function to produce an output. The weights determine the strength of the connections between the neurons and play a crucial role in the learning process.

The basic structure of a neural network consists of three fundamental components: input layer, hidden layer(s), and output layer. Each layer is composed of interconnected nodes called neurons or units. Here's a breakdown of the basic structure:

1.  Input Layer: The input layer is responsible for receiving the input data and passing it forward to the next layer. Each node in the input layer represents a feature or attribute of the input data. The number of nodes in the input layer corresponds to the dimensionality of the input data.

2.  Hidden Layer(s): Hidden layers are intermediate layers between the input and output layers. They perform computations on the input data and progressively extract higher-level features or representations. A neural network can have one or multiple hidden layers, depending on the complexity of the problem. Each neuron in a hidden layer receives inputs from the previous layer and applies a specific transformation, typically using an activation function.

3.  Output Layer: The output layer produces the final results or predictions based on the computations performed in the hidden layers. The number of nodes in the output layer depends on the nature of the problem being solved. For example, a binary classification task may have one output node representing the probability of belonging to one class, while a multi-class classification task may have multiple output nodes, each representing the probability of belonging to a specific class. Connections between neurons in different layers are represented by weights, which determine the strength or importance of the connection. Each connection has an associated weight that is multiplied by the input value at that connection. The weighted inputs are then summed up and passed through an activation function to introduce non-linearity into the network. Activation functions play a crucial role in introducing non-linear transformations within a neural network. They help the network learn complex patterns and make it capable of approximating non-linear functions. Common activation functions include the sigmoid function, ReLU (Rectified Linear Unit), tanh (hyperbolic tangent), and softmax (used for multi-class classification).The basic structure of a neural network can be visualized as a series of layers, with nodes in each layer connected to nodes in the adjacent layers. The weights associated with the connections are adjusted during the training process to minimize the error between the network's predictions and the desired output.

It's worth noting that this description covers the general structure of a feedforward neural network, which is the most common type. There are also other types of neural networks, such as recurrent neural networks (RNNs) and convolutional neural networks (CNNs), which have specific architectural variations suited for different tasks.Neurons are organized into layers within a neural network. The input layer receives the initial data, and the output layer produces the final results. In between the input and output layers, there can be one or more hidden layers. The hidden layers are responsible for extracting relevant features from the input data and learning complex patterns.The process of training a neural network involves adjusting the weights of the connections to minimize the difference between the predicted output and the desired output. This is typically done using an algorithm called backpropagation, which propagates the error backward through the network, updating the weights along the way.Neural networks have gained significant attention in recent years due to their ability to learn and solve complex problems, such as image recognition, natural language processing, and speech recognition. Deep learning, a subset of neural networks, involves training networks with multiple hidden layers and has achieved remarkable success in various fields.Neural networks play a crucial role in data mining and are widely used for various tasks within the field. Some ways neural networks are applied in data mining:

1) Pattern Recognition: Neural networks are used to identify and recognize patterns in large datasets. They can learn to recognize complex patterns and extract meaningful features from the input data, enabling tasks like image classification, speech recognition, and handwriting recognition.

2) Classification: Neural networks are used for classification tasks where the goal is to assign input data to predefined categories or classes. By training a neural network on labeled data, it can learn to classify new, unseen instances into the appropriate categories based on learned patterns and features.

3) Regression: Neural networks can be used for regression tasks, where the goal is to predict a continuous numeric value. By training a neural network on a dataset with input-output pairs, it can learn to approximate the underlying relationship between the inputs and outputs, allowing for the prediction of new output values for given inputs.

4) Clustering: Neural networks can also be employed for clustering tasks, which involve grouping similar instances together based on their characteristics or attributes. Self-organizing maps (SOM) are a type of neural network commonly used for clustering tasks. They can organize data into clusters based on their similarities and provide insights into the structure of the dataset.

5) Anomaly Detection: Neural networks are effective in detecting anomalies or outliers in data. By training a neural network on normal instances, it can learn to recognize patterns typical of the normal behavior and identify instances that deviate significantly from the learned patterns, indicating potential anomalies.

6) Association Rule Mining: Neural networks can be utilized to discover association rules, which reveal relationships and dependencies between variables or items in a dataset. By training a neural network to predict the occurrence of certain events or items based on the presence or absence of others, it can identify interesting associations and generate valuable insights.

Neural networks offer powerful modeling capabilities in data mining, allowing for complex pattern recognition, prediction, and exploration of large datasets. They can handle nonlinear relationships and adapt to changing patterns, making them versatile tools in extracting knowledge and making predictions from data.There are various neural network algorithms used in data mining, each designed to address different types of problems and data characteristics. Some commonly used neural network algorithms in data mining:

1) Multilayer Perceptron (MLP): MLP is a basic and widely used neural network algorithm. It consists of multiple layers of interconnected nodes (neurons) and uses backpropagation for training. MLP is suitable for tasks like classification and regression, and it can handle complex nonlinear relationships in data.A Multilayer Perceptron (MLP) is a type of artificial neural network that consists of multiple layers of nodes (neurons) in a directed graph, where each layer is fully connected to the next one. MLPs are the simplest form of deep learning models and are used primarily for supervised learning tasks such as classification and regression.The essential components of an MLP are:

1. **Input Layer:**
   - The input layer is the first layer of the network and consists of nodes that represent the features of the input data. Each node corresponds to one feature in the dataset.

2. **Hidden Layers:**
   - Hidden layers are intermediate layers between the input and output layers. Each hidden layer consists of multiple neurons, and each neuron in a layer is connected to every neuron in the previous layer (fully connected). The purpose of hidden layers is to capture complex patterns and interactions in the data.
   - MLPs can have one or more hidden layers, and the depth (number of hidden layers) and width (number of neurons in each layer) are hyperparameters that can be adjusted during model design.

3. **Output Layer:**
   - The output layer produces the final predictions of the network. For a classification task, the output layer typically has one neuron per class, and the output is a probability distribution

over the classes. For a regression task, the output layer usually has a single neuron that predicts a continuous value.

4. **Activation Functions:**

- Activation functions introduce non-linearity into the model, enabling the network to learn complex functions. Common activation functions used in MLPs include:

  - **ReLU (Rectified Linear Unit):** Outputs the input directly if it's positive; otherwise, it outputs zero. It's widely used in hidden layers due to its simplicity and effectiveness.

  - **Sigmoid:** Outputs a value between 0 and 1, often used in the output layer for binary classification tasks.

  - **Tanh:** Outputs values between -1 and 1, useful for zero-centered data.

5. **Weights and Biases:**

- **Weights:** Each connection between neurons has an associated weight, which determines the strength of the connection. During training, these weights are adjusted to minimize the error in predictions.

- **Bias:** Each neuron also has a bias term that allows the activation function to be shifted to fit the data better.

6. **Loss Function:**

- The loss function measures the difference between the network's predictions and the actual target values. Common loss functions include:

  - **Mean Squared Error (MSE):** Used for regression tasks.

  - **Cross-Entropy Loss:** Used for classification tasks, especially when the output is a probability distribution.

7. **Backpropagation:**

- Backpropagation is the algorithm used to train MLPs. It involves calculating the gradient of the loss function with respect to each weight in the network and using these gradients to update the weights in the opposite direction (gradient descent). This process is repeated iteratively to minimize the loss.

8. **Optimization Algorithms:**

- Optimization algorithms, such as Stochastic Gradient Descent (SGD), Adam, and RMSprop, are used to adjust the weights during training to minimize the loss function. Adam is particularly popular due to its adaptive learning rate and efficient convergence.

The working of an MLP is as follows:

1. **Forward Propagation:**
   - During forward propagation, the input data is passed through the network from the input layer to the output layer. Each neuron in a layer calculates a weighted sum of its inputs, applies an activation function, and passes the result to the next layer. The final output is generated at the output layer.

2. **Error Calculation:**
   - The error (difference between the predicted output and the actual target) is calculated using the loss function.

3. **Backpropagation and Weight Update:**
   - The error is propagated backward through the network to calculate the gradients of the loss with respect to each weight. The weights are then updated using the chosen optimization algorithm to reduce the error.

4. **Iteration:**
   - The process of forward propagation, error calculation, and backpropagation is repeated for a fixed number of iterations (epochs) or until the model converges (i.e., the loss stops decreasing).

**The various applications of MLPs are:**

1. **Classification:**
   - MLPs are commonly used for classification tasks, such as predicting whether an email is spam or not, or classifying handwritten digits. For example, an MLP with a softmax output layer can be used to classify images from the MNIST dataset.

2. **Regression:**
   - MLPs can also be used for regression tasks, where the goal is to predict a continuous value. For example, predicting house prices based on features like location, size, and age.

3. **Pattern Recognition:**
   - MLPs are used in pattern recognition tasks, such as recognizing speech patterns, detecting anomalies, or recognizing handwriting.

4. **Forecasting:**
   - MLPs can be applied to time-series forecasting tasks, where the goal is to predict future values based on past observations.

**The advantages and limitations of MLP are:**

**Advantages:**

- **Simplicity:** MLPs are relatively easy to understand and implement.
- **Versatility:** They can be applied to a wide range of supervised learning tasks.
- **Ability to Model Non-Linear Relationships:** MLPs can model complex, non-linear relationships between inputs and outputs.

**Limitations:**

- **Overfitting:** MLPs can easily overfit to training data, especially when the model is too complex or the dataset is small.
- **Computationally Expensive:** Training deep MLPs with many hidden layers can be computationally expensive and time-consuming.
- **Lack of Spatial Hierarchies:** Unlike Convolutional Neural Networks (CNNs), MLPs do not explicitly capture spatial hierarchies in data, making them less suitable for tasks like image recognition.

2) Convolutional Neural Networks (CNN): CNNs are primarily used for analyzing visual data, such as images and videos. They are designed to automatically learn and extract hierarchical patterns and features from the input data using convolutional layers. CNNs have been highly successful in image classification, object detection, and image recognition tasks.

Convolutional Neural Networks (CNNs) are a specialized type of neural network designed primarily for processing structured grid data, such as images. They are particularly effective in tasks involving image recognition, classification, object detection, and computer vision, but they can also be applied to other types of data like time series and audio. The key Concepts of CNNs are:

1. **Convolutional Layers:**
   - **Convolution Operation:** The core operation in a CNN, where a small matrix called a filter or kernel slides over the input data (e.g., an image) and performs element-wise multiplication and summation. This operation helps in capturing local patterns such as edges, textures, or color gradients.
   - **Filters/Kernels:** Learnable parameters that are trained to detect specific features in the input data. Different filters detect different features like edges, corners, and textures.

- **Feature Maps:** The output of the convolution operation, which is a transformed version of the input data highlighting the features detected by the filter.

2. **Pooling Layers:**
   - **Purpose:** Pooling layers reduce the dimensionality of the feature maps while retaining the most important information. This helps in reducing computational complexity and preventing overfitting.
   - **Max Pooling:** The most common type of pooling, where the maximum value within a small region (e.g., 2x2) of the feature map is selected and retained.
   - **Average Pooling:** Another pooling method where the average value within a small region is taken, though it is less commonly used than max pooling.

3. **Fully Connected Layers:**
   - After several convolutional and pooling layers, the high-level features are fed into one or more fully connected layers (dense layers). These layers act like traditional neural networks, where every neuron is connected to every neuron in the previous layer.
   - These layers are used to make the final prediction, whether it's classifying an image or detecting objects.

4. **Activation Functions:**
   - **ReLU (Rectified Linear Unit):** The most commonly used activation function in CNNs. It introduces non-linearity by outputting the input directly if it is positive, otherwise, it outputs zero.
   - **Softmax:** Often used in the final layer of a classification network to output a probability distribution over classes.

5. **Dropout:**
   - A regularization technique used in fully connected layers where a fraction of neurons are randomly turned off during training. This prevents the network from becoming too reliant on specific neurons and helps reduce overfitting.

6. **Stride and Padding:**
   - **Stride:** The number of pixels by which the filter moves across the input. A larger stride reduces the spatial dimensions of the output feature map.
   - **Padding:** Adding extra pixels (usually zeros) around the input to control the spatial size of the output feature map. Padding ensures that the spatial dimensions of the output are maintained or controlled.A typical CNN architecture consists of the following layers:

1. **Input Layer:**
   - Receives the raw image data, typically represented as a 3D array (height, width, and color channels).

2. **Convolutional Layer(s):**
   - Apply multiple filters to the input data, producing multiple feature maps. This layer captures the spatial hierarchies in the data, such as edges and textures.

3. **Pooling Layer(s):**
   - Reduces the spatial dimensions of the feature maps, retaining the most significant information and improving computational efficiency.

4. **Additional Convolutional and Pooling Layers:**
   - In deeper CNNs, multiple convolutional and pooling layers are stacked to capture increasingly complex features. Early layers may capture low-level features (edges), while deeper layers capture high-level features (object parts or entire objects).

5. **Fully Connected Layer(s):**
   - After the convolutional and pooling layers, the feature maps are flattened into a 1D vector and passed through fully connected layers. These layers combine the features learned to make a final prediction.

6. **Output Layer:**
   - The final fully connected layer, often followed by a softmax activation function, outputs the predicted class probabilities.

CNN has various applications such as:

1. **Image Classification:**
   - **Example:** Classifying images of animals into categories like cats, dogs, and birds. The ImageNet challenge popularized the use of CNNs for large-scale image classification tasks.

2. **Object Detection:**
   - **Example:** Identifying and localizing objects within an image, such as detecting cars and pedestrians in a self-driving car's camera feed. Techniques like YOLO (You Only Look Once) and R-CNN (Region-Based CNN) are used.

3. **Image Segmentation:**
   - **Example:** Classifying each pixel in an image to identify objects, such as segmenting a tumor in a medical image. Techniques like U-Net are commonly used.

4. **Face Recognition:**
   - **Example:** Identifying and verifying individuals' faces in photos or videos. CNNs are the backbone of many modern face recognition systems.

5. **Natural Language Processing (NLP):**
   - **Example:** Text classification, sentiment analysis, and language translation. Although more specialized architectures like RNNs and Transformers are often used, CNNs can be effective for certain NLP tasks.

6. **Audio Processing:**
   - **Example:** Speech recognition and music genre classification. CNNs can be applied to spectrograms (visual representations of audio signals) for processing audio data.Some popular CNN architectures are:

1. **LeNet-5:**
   - One of the earliest CNNs, designed for digit recognition in the MNIST dataset.

2. **AlexNet:**
   - A deep CNN that won the ImageNet competition in 2012, popularizing deep learning. It introduced the use of ReLU and dropout in CNNs.

3. **VGGNet:**
   - Known for its simplicity and depth, VGGNet uses very small (3x3) convolution filters and is deep, with up to 19 layers.

4. **GoogLeNet (Inception):**
   - Introduced the concept of Inception modules, which use multiple filter sizes at each layer, allowing the network to capture various levels of detail.

5. **ResNet (Residual Networks):**
   - Introduced the concept of residual connections (skip connections) to allow very deep networks (up to 152 layers) to be trained without suffering from the vanishing gradient problem.

6. **DenseNet:**
   - Connects each layer to every other layer in a feed-forward fashion, reducing the number of parameters and promoting feature reuse.

3) Recurrent Neural Networks (RNN): RNNs are suited for sequential data analysis, where the order of the data points matters, such as time series data or natural language processing tasks. RNNs

have recurrent connections that allow them to retain memory of previous inputs, making them effective in tasks like language modeling, speech recognition, and sentiment analysis.

4) Long Short-Term Memory (LSTM): LSTM is a variant of RNN that addresses the vanishing gradient problem and can effectively capture long-term dependencies in sequential data. LSTMs are widely used in tasks like machine translation, speech recognition, and text generation.

5) Radial Basis Function Networks (RBFN): RBFN is a type of neural network that uses radial basis functions as activation functions. It is often used for clustering and function approximation tasks. RBFNs have a hidden layer of radial basis functions that transform the input data into a higher-dimensional space, making them suitable for capturing complex relationships.

6) Self-Organizing Maps (SOM): SOM is an unsupervised learning algorithm used for clustering and visualization tasks. It uses a competitive learning process to create a low-dimensional representation of the input data, which helps reveal the underlying structure and relationships within the data.

**6.1.1 Backpropogation Algorithm**

Backpropagation, short for "backward propagation of errors," is an algorithm commonly used to train artificial neural networks. It enables the network to learn from labeled training data by adjusting the weights of the connections between neurons.The phases of backpropagation algorithm are:

1) Forward Propagation: The input data is presented to the neural network, and the activations of the neurons are computed layer by layer, starting from the input layer and progressing through the hidden layers to the output layer. The activations are calculated by applying the weighted sum of inputs to each neuron and passing it through an activation function.

2) Error Computation: The difference between the network's output and the expected output (target value) is calculated. This difference is known as the error or the loss. The specific error measure used depends on the task, such as mean squared error (MSE) for regression or cross-entropy for classification.

3) Backward Propagation: The error is propagated backward through the network, starting from the output layer towards the input layer. For each neuron, the error contribution is calculated by considering the weighted connections from the neurons in the subsequent layer and their associated errors. This is done using the chain rule of calculus.

4) Weight Update: The calculated error contributions are used to update the weights of the connections in the network. The weights are adjusted in the opposite direction of the gradient of the error with respect to the weights. The learning rate parameter determines the size of the weight

updates. A smaller learning rate results in smaller updates, while a larger learning rate can cause overshooting and instability.

5) Iteration: Steps 1 to 4 are repeated for multiple iterations or epochs, where the entire training dataset is presented to the network. This iterative process allows the network to gradually adjust the weights to minimize the error and improve its performance.

The backpropagation algorithm utilizes gradient descent optimization to iteratively update the weights in the direction that reduces the error. By iteratively adjusting the weights based on the computed errors, the network learns to approximate the desired output for a given input.The backpropagation is not limited to a single hidden layer. It can be applied to networks with multiple hidden layers, known as deep neural networks. In such cases, the algorithm propagates the errors backward through each layer, adjusting the weights layer by layer.Backpropagation has been a foundational algorithm in the field of neural networks and has contributed to the success of various deep learning architectures and applications

---

*Check your Progress 1*

1. What is a neural network? Describe its basic structure and functioning

2. What is the role of activation function?

3. Explain the process of forward propagation in a neural network.

4. Describe the backpropagation algorithm and its role in training neural networks.

4. List different types of layers in a neural network.

5. What is the role of activation function?

6. List the important application areas of artificial neural networks.

7. What is CNN?

8. Which algorithm is used to the neural networks?

---

# 6.3 Genetic Algorithms

Genetic algorithms are a class of optimization algorithms inspired by the principles of natural selection and genetics. They are commonly used in solving complex optimization problems where traditional algorithms may struggle to find the optimal solution.The main idea behind genetic algorithms is to simulate the process of evolution by iteratively improving a population of candidate solutions over successive generations. The algorithm starts with an initial population of potential solutions, which are represented as individuals or chromosomes. Each chromosome typically encodes a set of parameters or variables that define a potential solution to the problem.

The genetic algorithm operates through a series of steps:

1) Initialization: An initial population of chromosomes is randomly generated to start the process.

2) Evaluation: Each chromosome in the population is evaluated by a fitness function, which measures how well the solution performs in solving the problem. The fitness function assigns a numerical value, the fitness score, to each chromosome.

3) Selection: Chromosomes with higher fitness scores are more likely to be selected for reproduction. Selection methods like tournament selection or roulette wheel selection are used to choose the fittest individuals.

4) Crossover: The selected chromosomes undergo crossover or recombination. This involves exchanging genetic information between pairs of chromosomes to create new offspring. The crossover can occur at a specific point or be done probabilistically across the chromosomes.

5) Mutation: To introduce diversity and explore new regions of the solution space, random changes or mutations are applied to some individuals in the population. These changes can alter certain genes or parameters within a chromosome.

6) Replacement: The offspring, which includes the crossover results and mutated individuals, replace a portion of the previous generation. This ensures that the population evolves over time towards better solutions.

7) Termination: The algorithm continues iterating through the selection, crossover, and mutation steps until a termination condition is met. This condition can be a maximum number of generations, reaching a desired fitness threshold, or a specific runtime limit.By iteratively applying these steps, genetic algorithms explore the search space and converge towards optimal or near-optimal solutions. They are particularly useful when dealing with large solution spaces, combinatorial optimization problems, or situations where the problem does not have a well-defined mathematical form.

Genetic algorithms have been successfully applied in various domains, including engineering design, scheduling, financial modeling, and artificial intelligence. They offer a flexible and robust approach to optimization and can handle complex problems that are difficult to solve using traditional methods.Genetic algorithms can be applied to various data mining tasks, particularly in optimization problems. Some ways genetic algorithms are used in data mining:

1) Feature Selection: Feature selection aims to identify the most relevant features from a large set of variables or attributes. Genetic algorithms can be used to search through the space of possible feature subsets and determine the optimal combination of features that maximizes the performance of a data

mining model. The genetic algorithm evaluates different subsets based on a fitness function, which can be derived from the performance of a classifier or a regression model.

2) Clustering: Genetic algorithms can be employed to optimize clustering algorithms by determining the best configuration of parameters. The genetic algorithm can search for the optimal number of clusters, the initial centroids, or the distance metric used in the clustering process. By iteratively evaluating different parameter combinations, the genetic algorithm can enhance the effectiveness of clustering algorithms.

3) Association Rule Mining: Association rule mining aims to discover relationships or associations between items or attributes in a dataset. Genetic algorithms can be used to optimize the rule discovery process by adjusting parameters such as the minimum support and confidence thresholds. The genetic algorithm searches for the optimal parameter settings that yield the most interesting and meaningful association rules.

4) Neural Network Architecture Optimization: Genetic algorithms can be employed to optimize the architecture of neural networks for data mining tasks. The genetic algorithm can explore various combinations of network topologies, activation functions, and learning parameters to find the best configuration that maximizes the neural network's performance on a given dataset.

5) Data Preprocessing: Genetic algorithms can assist in optimizing data preprocessing steps such as data cleaning, normalization, and feature scaling. The genetic algorithm can search for the best combination of preprocessing operations and parameter settings that lead to improved data quality and enhanced performance of subsequent data mining algorithms.

By using genetic algorithms in data mining, researchers and practitioners can automate the process of searching for optimal configurations or solutions, particularly in cases where the search space is large and complex. Genetic algorithms provide an efficient and effective approach for finding near-optimal solutions and exploring different possibilities in the data mining process.

## 6.4 ROUGH SETS

Rough sets are a mathematical framework used for dealing with uncertainty and imprecision in data analysis and decision-making. The concept of rough sets was introduced by Zdzisław Pawlak in the early 1980s as an extension of classical set theory.

The main idea behind rough sets is to approximate and describe the boundaries or discernibility regions of different classes or concepts within a dataset. It allows for the analysis of data when there is incomplete or uncertain information available. Rough set theory provides a formal mathematical framework to handle

this uncertainty and make data-driven decisions.In rough set theory, a dataset is represented as a collection of objects or instances, each characterized by a set of attributes or features. The attributes can be binary or categorical, and they define the properties or characteristics of the objects. Rough sets aim to identify subsets of attributes that are sufficient for distinguishing between different classes or concepts. The key concepts in rough set theory include:

1.      Lower Approximation: The lower approximation of a class is the set of objects for which we can be certain that they belong to the class based on the available data.

2.      Upper Approximation: The upper approximation of a class is the set of objects for which there is evidence or support to consider them as belonging to the class, but it is not certain.

3.      Boundary Region: The boundary region consists of the objects that are on the edge or borderline between two or more classes. These objects exhibit uncertain or ambiguous characteristics and cannot be definitively assigned to a single class.

4.      Decision Rules: Decision rules are derived from the lower and upper approximations and define the conditions or attribute combinations that can be used to classify objects into different classes. Decision rules capture the patterns and dependencies present in the data and provide a basis for making decisions or predictions.

Rough set theory has been applied in various domains, including data mining, machine learning, pattern recognition, and expert systems. It provides a useful framework for feature selection, attribute reduction, rule induction, and knowledge discovery from data. By capturing the uncertainty and imprecision in data, rough sets help in understanding the structure and relationships within complex datasets and support decision-making processes in real-world applications.In data mining, rough set theory can be applied to various tasks and provide valuable insights into the structure of complex datasets. Here are some specific applications of rough sets in data mining:

1.      Feature Selection: Rough set theory can be used to identify relevant features or attributes from a large dataset. By analyzing the lower and upper approximations of different classes, rough sets can help determine which attributes are essential for distinguishing between classes and can be used for effective feature selection.

2.      Attribute Reduction: Rough sets can assist in reducing the dimensionality of a dataset by eliminating redundant or irrelevant attributes. By examining the dependencies and relationships between attributes, rough set-based attribute reduction techniques can identify subsets of attributes that retain the discriminative power of the dataset while reducing its complexity.

3.      Rule Induction: Rough sets provide a framework for inducing decision rules from data. By analyzing the lower and upper approximations of different classes, rough set-based rule induction algorithms can identify conditions or attribute combinations that lead to accurate classification or prediction. These decision rules capture the patterns and dependencies present in the data and can be used for knowledge discovery and decision-making.

4.      Data Clustering: Rough set theory can be employed for clustering analysis to identify groups or clusters of similar objects. By examining the lower and upper approximations of different clusters, rough set-based clustering algorithms can determine the boundaries and characteristics of each cluster, providing insights into the underlying structure of the dataset.

5.      Association Rule Mining: Rough sets can be used to discover association rules that reveal relationships between attributes or items in a dataset. By analyzing the dependencies and co-occurrence patterns in the lower and upper approximations, rough set-based association rule mining techniques can identify interesting and meaningful associations among attributes or items.

6.      Missing Data Handling: Rough set theory provides a framework for dealing with missing or incomplete data. By analyzing the lower and upper approximations of objects with missing values, rough sets can estimate the potential values or fill in the missing data based on the available information, allowing for more robust analysis and modeling. Rough set theory offers a flexible and interpretable approach to data mining, particularly in situations where there is uncertainty, imprecision, or incomplete information. It helps in understanding the structure and dependencies within complex datasets, facilitating knowledge discovery and decision-making processes.

---

*Check your progress 2*

Q1. Explain the following terms in genetic algorithms: gene, chromosome, and fitness.

Q2. What is the basic idea behind genetic algorithm?

Q3. Why crossover operation is performed?

Q4. List the major applications of genetic algorithms.

Q5. What is the application of rough sets?

---

## 6.5  Summary

In this unit, the concepts of Neural network, Genetic algorithms and rough sets was explained. The structure of neural network and its application in data mining is coveredin  this unit. Back propagation

algorithm is also explained in this unit. The concept of genetic algorithms and their various functions are also discussed in this unit.

The unit explained the following:

- Neural networks
- Genetic algorithms
- Rough Sets

**Suggested Further Reading**

1. "Data Mining: Concepts and Techniques" by Jiawei Han, Micheline Kamber, and Jian Pei: This book provides a comprehensive introduction to the concepts, methods, and algorithms of data mining. It covers various data mining techniques, including classification, clustering, association analysis, and outlier detection.

2. "Introduction to Data Mining" by Pang-Ning Tan, Michael Steinbach, and Vipin Kumar: This book provides a comprehensive introduction to the fundamental concepts and techniques of data mining. It covers topics such as data preprocessing, classification, clustering, association analysis, and anomaly detection.

3. "Data Mining: Concepts, Models, Methods, and Algorithms" by Mehmed Kantardzic: This book presents a comprehensive overview of data mining concepts, models, methods, and algorithms. It covers key topics such as data preprocessing, classification, clustering, pattern mining, and data visualization.

## Review Questions

Q1. What is neural network? Discuss Back-propagation algorithm in neural network.

_____

_____

_____

Q2. What are genetic algorithms? Discuss the applicationof genetic algorithms in data mining process.

_____

_____

_____

Q3. Explain various applications of rough sets in data mining.

_____

_____

_____

Q4. Discuss various applications of classification in data mining.

_____

_____

**Master of Computer Application**

# MCA-E6N
# Data Mining

उ० प्र० राजर्षि टण्डन
मुक्त विश्वविद्यालय, प्रयागराज

# Block

# 3

## Data-Mining Techniques in Detail

## Course Design Committee

| | |
|---|---|
| **Prof. Ashutosh Gupta** | **Chairman** |
| Director (In-charge) | |
| School of Computer & Information Science, UPRTOU Allahabad | |
| **Prof. Suneeta Agarwal** | **Member** |
| Dept. of Computer Science & Engineering | |
| Motilal Nehru National Institute of Technology Allahabad | |
| **Dr. Upendra Nath Tripathi** | **Member** |
| Associate Professor | |
| DeenDayalUpadhyay Gorakhpur University, Gorakhpur | |
| **Dr. Ashish Khare** | **Member** |
| Associate Professor | |
| Dept. of Computer Science, University of Allahabad, Prayagraj | |
| **Ms. Marisha** | **Member** |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |
| **Mr. Manoj Kumar Balwant** | **Member** |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |

## Course Preparation Committee

| | |
|---|---|
| **Dr. Tulika Narang** | **(Block 1, 2 & 4) Author** |
| Assistant Professor, Computer science | |
| United University, Rawatpur, Prayagraj | |
| **Dr. Krishan Kumar** | **(Block 3-Unit 7) Author** |
| Assistant Professor, | |
| Department of Computer Science, Faculty of Technology | |
| Gurukula Kangri Vishwavidyalaya, Haridwar (UK) | |
| **Dr. Pooja Yadav** | **(Block 3-Unit 8) Author** |
| Assistant, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Iram Naim** | **(Block 3-Unit 9) Author** |
| Assistant Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Brajesh Kumar** | **Editor** |
| Associate Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Mr. Manoj Kumar Balwant** | **Coordinator** |
| Assistant Professor (Computer Science), | |
| School of Sciences, UPRTOU, Prayagraj | |

This block's first unit explores foundational concepts of descriptive analytics and cluster analysis, critical for interpreting data patterns and making informed decisions. Descriptive analytics involves analyzing historical data to summarize past performance, identify trends, and detect anomalies through data aggregation, statistical measures, and visualization techniques. This transforms raw data into actionable insights, typically presented via dashboards and reports.

Cluster analysis, a statistical method for grouping similar data points, aids in discovering patterns without predefined labels. Various algorithms, such as K-means, hierarchical clustering, and DBSCAN, identify data structures with market segmentation and anomaly detection applications. Partitioning-based algorithms minimize variance within clusters, while hierarchical methods illustrate relationships through tree structures. Density-based approaches, like DBSCAN, effectively identify clusters of varying shapes based on data density.

The eighth unit shifts focus to supervised learning algorithms for classification and prediction. It covers decision trees, lazy and eager learners, Bayesian classification, and rule-based models. Additionally, it discusses backpropagation in neural networks and the application of Support Vector Machines (SVM) for regression and classification tasks.

Finally, the last unit addresses algorithms for mining frequent item sets, each with unique strengths and weaknesses. Techniques such as the Apriori and FP-Growth algorithms, Eclat, Direct Hashing and Pruning (DHP), and Rapid Association Rule Mining (RARM) facilitate the discovery of complex patterns in various fields, enhancing decision-making through context-aware insights. Understanding the distinctions between association rule mining and correlation analysis is essential for practical data analysis and applying the appropriate techniques.

# Unit 7: Descriptive Analytics - Cluster Analysis

**Structure**

## 7.0 Introduction To Descriptive Analytics

Descriptive analytics is a fundamental aspect of data mining that involves summarizing and interpreting historical data to uncover patterns, trends, and insights. Unlike predictive or prescriptive analytics, which aim to forecast future events or recommend actions, descriptive analytics focuses on understanding what has happened in the past. This type of analysis helps organizations make sense of large datasets by providing a clear, concise overview of the information contained within. It forms the foundation for more advanced data analysis techniques and supports decision-making by providing a solid understanding of historical data.

### 7.0.1 Key Techniques in Descriptive Analytics

*7.0.1.1 Summary Statistics*

- **Mean, Median, and Mode:** These measures of central tendency provide insights into the typical values within a dataset. The mean gives the average, the median provides the middle value, and the mode indicates the most frequently occurring value.

- **Standard Deviation and Variance:** These measures of dispersion describe the spread of data points around the mean. A high standard deviation indicates that data points are spread out over a wide range of values, while a low standard deviation indicates that they are clustered closely around the mean.

### 7.0.1.2 Data Visualization

- **Histograms:** These graphical representations show the distribution of a single variable, helping to identify patterns such as skewness and modality.

- **Bar Charts and Pie Charts:** These are useful for displaying categorical data, making it easy to compare different categories and see their relative frequencies.

- **Line Graphs:** These are ideal for illustrating trends over time, allowing for easy identification of upward or downward trends.

- **Scatter Plots:** These plots show the relationship between two variables, helping to identify correlations and potential causations.

- **Box Plots:** These visualizations display the distribution of data based on five summary statistics: minimum, first quartile, median, third quartile, and maximum, highlighting outliers and the spread of the data.

### 7.0.1.3 Frequency Distribution

Frequency distribution involves counting the number of times each value occurs in a dataset. This technique is particularly useful for categorical data and can be represented in tables or charts. It helps to identify the most and least common values and to understand the distribution of data.

### 7.0.1.4 Cross-Tabulation

This is also known as contingency tables; cross-tabulation examines the relationship between two or more categorical variables. It displays the frequency distribution of variables in a matrix format, allowing for easy comparison and identification of potential correlations. This technique is particularly useful in market research and survey analysis.

### 7.0.1.5 Clustering and Segmentation

Clustering techniques like *k-means* clustering group similar data points together based on their attributes. This helps in identifying distinct segments within the data, which can be further analyzed for targeted marketing strategies or personalized services. While clustering is often associated with more advanced analytics, its basic application in descriptive analytics provides valuable insights into the structure of the dataset.

### 7.0.1.6 Correlation Analysis

Correlation analysis measures the strength and direction of the relationship between two variables. The correlation coefficient, ranging from -1 to 1, indicates whether variables move together in the same direction (positive correlation) or in opposite directions (negative correlation). This technique helps in understanding the degree to which variables are related, which can be crucial for identifying factors that influence outcomes.

## 7.0.2 Applications of Descriptive Analytics

### 7.0.2.1 Business and Marketing

Descriptive analytics helps companies understand customer behavior, preferences, and purchasing patterns. By analyzing past sales data, businesses can identify which products are most popular, which customer segments are most profitable, and which marketing campaigns have been most effective. This information is crucial for developing targeted marketing strategies and optimizing product offerings.

### 7.0.2.2 Healthcare

In healthcare, descriptive analytics is used to understand patient demographics, disease prevalence, and treatment outcomes. Hospitals can analyze patient records to determine the most common illnesses, track the effectiveness of different treatments, and identify areas for improvement in patient care. This data-driven approach helps in improving healthcare services and patient outcomes.

### 7.0.2.3 Finance

Financial institutions use descriptive analytics to monitor market trends, assess risk, and track investment performance. By analyzing historical stock prices, transaction data, and economic indicators, analysts can gain insights into market conditions and make informed investment decisions. Descriptive analytics also helps in identifying fraudulent activities and managing financial risks.

### 7.0.2.4 Education

Educational institutions employ descriptive analytics to evaluate student performance, enrollment trends, and program effectiveness. This analysis helps in identifying areas where students struggle and developing targeted interventions to improve educational outcomes. It also aids in optimizing resource allocation and enhancing the overall quality of education.

### 7.0.2.5 Public Policy

Governments and policymakers use descriptive analytics to understand social and economic conditions, such as unemployment rates, crime statistics, and population demographics. This information is critical for developing policies, allocating resources effectively, and addressing societal issues. Descriptive analytics helps in making data-driven decisions that can positively impact communities.

### 7.0.2 Challenges and Considerations

*7.0.2.1 Data Quality*

The accuracy and reliability of descriptive analytics depend heavily on the quality of the data. Inaccurate, incomplete, or biased data can lead to misleading conclusions. Ensuring data quality through proper cleaning and validation processes is essential to produce meaningful and trustworthy results.

*7.0.2.2 Contextual Understanding*

Descriptive analytics provides a snapshot of the data, but interpreting the results requires contextual knowledge. Understanding the underlying factors that drive the observed patterns is crucial for making meaningful inferences. Analysts must consider the broader context and domain-specific knowledge to draw accurate conclusions from the data.

*7.0.2.3 Over-Simplification*

There is a risk of oversimplifying complex data when relying solely on summary statistics or basic visualizations. Important nuances and details can be overlooked, leading to incomplete or incorrect interpretations. Combining multiple techniques and considering various perspectives can help mitigate this risk.

*7.0.2.4 Dynamic Data*

In rapidly changing environments, such as financial markets or social media, descriptive analytics based on static historical data may quickly become outdated. Continuous monitoring and updating of data are necessary to maintain relevance and ensure that the insights derived from the analysis remain applicable

Consequently, it can be said that Descriptive analytics is a foundational element of data mining, providing essential insights into the structure and characteristics of data. By employing techniques such as summary statistics, data visualization, frequency distribution, cross-tabulation, clustering, and correlation analysis, it enables businesses and researchers to understand and interpret complex datasets. Despite its challenges, when executed effectively, descriptive analytics serves as a critical first step in the data mining process, paving the way for more advanced analytical approaches and informed decision-making. Through its applications in various industries, descriptive analytics helps organizations make sense of their data, uncover valuable insights, and drive strategic actions.

## 7.1 Objectives

At the end of this unit you will come to know about the following:

- Introduction of descriptive analytics

- Types of clustering algorithms

- Partitioning-based algorithm

- Density-based algorithm

- Grid-based algorithm

- Hierarchical-based algorithm

- Model based algorithm

- Outlier analysis

## 7.2 Introduction To Clustering Analysis

Clustering analysis is a technique in data mining and machine learning that groups a set of objects in such a way that objects in the same group (called a cluster) are more similar to each other than to those in other groups (clusters). It is a form of unsupervised learning, meaning that it deals with unlabeled data, and it is used to find natural groupings or patterns within data. Clustering analysis has a wide range of applications across various fields. In marketing, it is used to segment customers based on purchasing behavior, enabling targeted marketing strategies. In biology, it helps in classifying organisms into species or grouping genes with similar expression patterns. In finance, clustering can detect fraudulent transactions by identifying patterns that deviate from normal behavior. In the field of social networks, it is used to identify communities within large networks.

## 7.2.1 Types of Clustering Methods

### 7.2.1 Partitioning Methods

- **K-means Clustering:** This is one of the most popular clustering algorithms. It aims to partition the dataset into $k$ clusters, where each cluster is represented by the mean of its points. The algorithm iteratively assigns data points to the nearest cluster center and then updates the cluster centers based on the points assigned to them.

- **K-medoids Clustering:** Similar to K-means, but instead of using the mean of the points, it uses an actual data point as the center of the cluster (medoid). This method is more robust to noise and outliers.

### 7.2.2 Hierarchical Methods

- **Agglomerative Clustering:** This is a bottom-up approach where each data point starts in its own cluster, and pairs of clusters are merged as one moves up the hierarchy. The merging continues until all points are in a single cluster or until a certain number of clusters is reached.

- **Divisive Clustering:** This is a top-down approach, which starts with all data points in a single cluster and recursively splits the clusters into smaller clusters until each point is in its own cluster or a predefined number of clusters is achieved.

- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** This method defines clusters as areas of high density separated by areas of low density. It can find arbitrarily shaped clusters and can identify noise or outliers.

- **OPTICS (Ordering Points to Identify the Clustering Structure):** An extension of DBSCAN, OPTICS works well with varying densities and can produce a more detailed clustering structure.

*7.2.4. Model-Based Methods*

- **Gaussian Mixture Models (GMMs):** These assume that the data is generated from a mixture of several Gaussian distributions. The Expectation-Maximization (EM) algorithm is typically used to find the parameters of these distributions.

- **Bayesian Clustering:** This method incorporates prior knowledge or beliefs about the data distribution and updates these beliefs as more data becomes available.

## 7.2.2 Evaluating Clustering Quality

Evaluating the quality of clustering results is challenging because it is an unsupervised method. Common evaluation metrics include:

- **Silhouette Score:** Measures how similar an object is to its own cluster compared to other clusters. A higher silhouette score indicates better-defined clusters.

- **Davies-Bouldin Index:** Computes the average similarity ratio of each cluster with the one most similar to it. Lower values indicate better clustering.

- **Dunn Index:** Measures the ratio of the minimum inter-cluster distance to the maximum intra-cluster distance. Higher values suggest well-separated clusters.

- **Elbow Method:** Used with K-means clustering to determine the optimal number of clusters by plotting the total within-cluster sum of square (WCSS) against the number of clusters and looking for an "elbow" point where the WCSS decreases abruptly.

## 7.2.3 Challenges and Considerations

Several challenges arise in clustering analysis, including:

- **Determining the number of clusters:** Many algorithms require the number of clusters to be specified in advance, which may not always be known.

- **Handling high-dimensional data:** High-dimensional data can lead to the curse of dimensionality, where the distance measures become less meaningful.

- **Scalability:** For very large datasets, the computational cost of clustering can be high.

- **Interpretability:** Understanding and interpreting the clusters can be complex, especially for high-dimensional data.

Clustering analysis is a powerful tool for discovering patterns and groupings in data. Its applications span multiple domains, from market segmentation to bioinformatics. Despite its challenges, advancements in algorithms and computational power continue to enhance its effectiveness and applicability. Understanding different clustering methods and evaluation metrics is crucial for selecting the appropriate technique for specific data and objectives.

## 7.3 Clustering Algorithms

Clustering algorithms are techniques used to group data points into clusters, where points within the same cluster are more similar to each other than to those in other clusters. These algorithms can be broadly categorized into several types: partitioning, hierarchical, density-based, and model-based. Partitioning algorithms, like K-Means and K-Medoids, divide the data into a predefined number of clusters by optimizing an objective function. Hierarchical algorithms, such as Agglomerative and Divisive clustering, build a hierarchy of clusters that can be visualized with a dendrogram. Density-based algorithms, like DBSCAN and OPTICS, identify clusters based on the density of data points, effectively handling noise and discovering clusters of varying shapes. Model-based algorithms, such as Gaussian Mixture Models, assume a probabilistic model for the data and estimate parameters to find the best fit. Each type of clustering algorithm has its strengths and weaknesses, and the choice of algorithm depends on factors such as the nature of the data, the desired cluster characteristics, and computational efficiency.

### 7.3.1 K-Means Clustering

- **Description:** Partitions the data into $K$ clusters by minimizing the sum of squared distances between data points and the centroids of their assigned clusters.

- **Pros:** Simple and efficient for large datasets.

- **Cons:** Requires the number of clusters $K$ to be specified in advance; sensitive to initial placement of centroids; assumes spherical clusters.

### 7.3.2 Hierarchical Clustering

- **Description:** Builds a tree (or dendrogram) of clusters by either a bottom-up approach (agglomerative) or a top-down approach (divisive).

- **Pros:** Does not require the number of clusters to be specified; dendrogram can be cut at different levels for different numbers of clusters.

- **Cons:** Computationally expensive for large datasets; difficult to scale.

### 7.3.3 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

- **Description:** Groups together points that are closely packed together, marking points in low-density regions as outliers.

- **Pros:** Can find clusters of arbitrary shape; does not require the number of clusters to be specified in advance.

- **Cons:** Requires careful selection of parameters (eps and minPts); struggles with varying density.

## 7.3.4 Mean Shift Clustering

- **Description:** Iteratively shifts data points towards the mode (highest density point) to find clusters.

- **Pros:** Does not require specifying the number of clusters; can find clusters of arbitrary shape.

- **Cons:** Computationally intensive; sensitive to bandwidth parameter.

## 7.3.5 Gaussian Mixture Models (GMM)

- **Description:** Assumes that the data is generated from a mixture of several Gaussian distributions with unknown parameters.

- **Pros:** Can model clusters with different shapes; probabilistic framework provides soft clustering.

- **Cons:** Requires the number of clusters to be specified in advance; can be sensitive to initial parameters.

## 7.3.6 Agglomerative Clustering

- **Description:** A type of hierarchical clustering that starts with each data point as a single cluster and merges the closest pairs of clusters iteratively.

- **Pros:** Easy to implement; produces a hierarchical tree.

- **Cons:** Computationally expensive for large datasets; requires a method to determine the optimal number of clusters.

## 7.3.7 BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies)

- **Description:** Constructs a tree structure (CF tree) to incrementally and dynamically cluster incoming data points.

- **Pros:** Efficient for large datasets; handles noise well.

- **Cons:** Requires a good initial choice of parameters; less effective for small datasets.

## 7.3.8 Spectral Clustering

- **Description:** Uses eigenvalues of a similarity matrix to perform dimensionality reduction before clustering in fewer dimensions.

- **Pros:** Can capture complex relationships between data points; effective for non-convex clusters.

- **Cons:** Computationally expensive; requires specifying the number of clusters.

## 7.3.9 OPTICS (Ordering Points to Identify the Clustering Structure)

- **Description:** Similar to DBSCAN but can identify clusters with varying densities.

- **Pros:** Can handle varying densities; no need to specify the number of clusters.

- **Cons:** More complex and slower than DBSCAN; sensitive to parameter settings.

## 7.3.10 Applications

Clustering algorithms are used in various fields, including:

- **Market Segmentation:** Grouping customers based on purchasing behavior.

- **Image Segmentation:** Partitioning images into regions of interest.

- **Anomaly Detection:** Identifying outliers or unusual data points.

- **Document Clustering:** Organizing large collections of text documents.

- **Biological Data Analysis:** Grouping genes or proteins with similar expression patterns.

Each algorithm has its strengths and weaknesses, making some more suitable for certain types of data and problems than others. Choosing the right algorithm often involves experimenting with different methods and tuning their parameters to fit the specific characteristics of the dataset.

# 7.4 Partitioning Based Algorithm

Partitioning-based clustering algorithms aim to divide a dataset into a set of distinct, non-overlapping clusters. The goal is to assign each data point to exactly one cluster, such that the data points within a cluster are more similar to each other than to those in other clusters. The most well-known partitioning algorithm is K-Means, but there are other methods as well.

## 7.4.1 K-Means Clustering

K-Means clustering is a popular algorithm used to partition data into a specified number of clusters, k, by iteratively minimizing the variance within each cluster. It begins by selecting k initial centroids and assigns each data point to the nearest centroid based on a distance metric, typically Euclidean distance. After assigning points, the centroids are recalculated as the mean of the points in each cluster, and the process repeats until convergence, when the centroids stabilize or a set number of iterations is reached. While K-Means is valued for its simplicity and efficiency, it requires the number of clusters to be predetermined and can be sensitive to initial centroid positions, potentially leading to suboptimal results. It also assumes clusters are spherical and of similar size, which may not always align with the true structure of the data.

This include the following steps:

- **Initialization:** Select $K$ initial centroids randomly or using some heuristic.

- **Assignment:** Assign each data point to the nearest centroid, forming *K* clusters.

- **Update:** Recalculate the centroids as the mean of all data points in each cluster.

- **Repeat:** Iterate the assignment and update steps until convergence (i.e., when assignments no longer change or the change is below a certain threshold).

The pseudocode for the K-Means clustering can be given as:

```
def k_means(X, K, max_iters=100):

    # Step 1: Initialize centroids
    centroids = initialize_centroids(X, K)

    for i in range(max_iters):

        # Step 2: Assign clusters
        clusters = assign_clusters(X, centroids)

        # Step 3: Update centroids
        new_centroids = update_centroids(X, clusters, K)

        # Check for convergence
        if has_converged(centroids, new_centroids):
            break

        centroids = new_centroids

    return clusters, centroids
```

As fasr as the pros of K-Means clustering are concerened; it is quite simple and easy to implement and efficient for large datasets. Moreover, it works well when clusters are spherical and of roughly equal size. On the other hand, if we talk about the cons, it requires the number of clusters *K* to be specified in advance and sensitive to the initial placement of centroids. Clusters should be spherical and equally sized. Furthermore, it is not suitable for clusters with varying densities or non-convex shapes. However, this is most commonly used clustering algorithm in machine learning applications.

## 7.4.2 K-Medoids Clustering

K-Medoids clustering is a robust method for partitioning a dataset into a specified number of clusters by selecting actual data points as cluster centers, known as medoids. Unlike k-means, which uses the mean of points to represent a cluster center, k-medoids chooses real data points that minimize the total dissimilarity within each cluster. This approach makes k-medoids less sensitive to outliers and noise, as it is based on actual data points rather than the average. The algorithm typically iterates between assigning data points to the nearest medoid and updating the medoids to minimize the sum of dissimilarities. While

k-medoids can handle non-spherical clusters and outliers better than k-means, it can be computationally expensive for large datasets due to its higher complexity in medoid selection and distance calculations. Despite this, k-medoids is valuable in applications where robustness to outliers and data integrity are critical.

### 7.4.2.1 Steps

- **Initialization:** Select $K$ initial medoids randomly.

- **Assignment:** Assign each data point to the nearest medoid.

- **Update:** For each medoid, consider swapping it with a non-medoid point and calculate the total cost of the clustering (sum of distances). Choose the swap that results in the lowest cost.

- **Repeat:** Iterate the assignment and update steps until convergence.

### 7.4.2.2 Algorithm

```
def k_medoids(X, K, max_iters=100):

    # Step 1: Initialize medoids
    medoids = initialize_medoids(X, K)

    for i in range(max_iters):
        # Step 2: Assign clusters
        clusters = assign_clusters(X, medoids)

        # Step 3: Update medoids
        new_medoids = update_medoids(X, clusters, medoids)

        # Check for convergence
        if has_converged(medoids, new_medoids):
            break

        medoids = new_medoids
    return clusters, medoids
```

### 7.4.2.3 Pros

- More robust to noise and outliers compared to K-Means.

- Can handle clusters with different shapes and sizes better than K-Means.

### 7.4.2.4 Cons

- More computationally expensive than K-Means.

- Still requires the number of clusters K to be specified in advance.

## 7.4.3 CLARA (Clustering Large Applications)

CLARA (Clustering Large Applications) is an extension of K-Medoids designed to handle large datasets. It works by applying K-Medoids to multiple samples of the dataset to find a good set of medoids.

### 7.4.3.1 Steps

- **Sampling:** Select multiple random samples from the dataset.

- **Apply K-Medoids:** Run K-Medoids on each sample to identify a set of medoids.

- **Evaluate:** Evaluate the medoids on the entire dataset and select the best set of medoids.

- **Assignment:** Assign all data points to the nearest medoid.

### 7.4.3.2 Algorithm

```python
def clara(X, K, num_samples=5, sample_size=40, max_iters=100):
    best_medoids = None
    best_cost = float('inf')

    for _ in range(num_samples):
        sample = random_sample(X, sample_size)
        medoids = k_medoids(sample, K, max_iters)[1]

        # Evaluate on the entire dataset
        clusters = assign_clusters(X, medoids)
        cost = calculate_total_cost(X, clusters, medoids)

        if cost < best_cost:
            best_cost = cost
            best_medoids = medoids

    clusters = assign_clusters(X, best_medoids)
    return clusters, best_medoids
```

### 7.4.3.3 Pros

- Scales better to large datasets than K-Medoids.

- Combines the robustness of K-Medoids with sampling to improve efficiency.

### 7.4.3.4 Cons

- Still requires the number of clusters K to be specified in advance.

- Performance depends on the quality and representativeness of the samples.

## 7.4.4 Applications

Partitioning-based clustering algorithms are used in various applications, including:

- **Market Segmentation:** Identifying distinct groups of customers for targeted marketing.

- **Image Segmentation:** Dividing images into regions for analysis.

- **Document Clustering:** Organizing large sets of documents into thematic clusters.

- **Anomaly Detection:** Identifying unusual patterns in data.

These algorithms are fundamental tools in data mining, providing efficient and intuitive ways to uncover structure in complex datasets.

---

## Check Your Progress

- Write the significance of cluster analysis in data mining applications.

- What do you understand by a cluster?

- Give the names of algorithms associated with partition based approach.

---

# 7.5 Hierarchical Based

Hierarchical-based algorithms are clustering techniques that build a hierarchy of clusters either through a divisive (top-down) or agglomerative (bottom-up) approach. In agglomerative clustering, each data point starts as its own cluster, and pairs of clusters are iteratively merged based on similarity until a single cluster remains or a stopping criterion is met. Conversely, divisive clustering begins with all data points in one cluster and recursively splits it until each point is isolated or a predefined condition is reached. These algorithms generate a dendrogram, a tree-like diagram that illustrates the merging or splitting process, providing a visual representation of the data's hierarchical structure. Hierarchical-based methods are valuable for their interpretability and ability to reveal nested cluster structures, but they can be computationally expensive and less scalable to large datasets. They are particularly useful in exploratory data analysis, where understanding the relationships between clusters at different levels of granularity is crucial.

## 7.5.1 Agglomerative Hierarchical Clustering

Agglomerative clustering starts with each data point as a separate cluster and then merges the closest pairs of clusters iteratively until all points are in a single cluster or a stopping criterion is met.

*7.5.1.1 Steps*

- **Initialization:** Start with each data point as a singleton cluster.

- **Compute Proximity:** Calculate the proximity matrix for all clusters (typically using a distance metric like Euclidean distance).

- **Merge Clusters:** Find the pair of clusters with the smallest distance and merge them.

- **Update Proximity Matrix:** Recalculate the proximity matrix to reflect the merge.

- **Repeat:** Continue merging until all data points are in one cluster or the desired number of clusters is achieved.

*7.5.1.2 Common Linkage Methods*

- **Single Linkage (Minimum Linkage):** Distance between the closest points of the clusters.

- **Complete Linkage (Maximum Linkage):** Distance between the farthest points of the clusters.

- **Average Linkage:** Average distance between all pairs of points in the two clusters.

- **Ward's Method:** Minimizes the total within-cluster variance.

## 7.5.2 Divisive Hierarchical Clustering

Divisive clustering, also known as "top-down" clustering, starts with all data points in a single cluster and recursively splits the clusters until each point is in its own cluster or another stopping criterion is met. This method is the opposite of agglomerative clustering, which starts with individual points and merges them. Divisive clustering can be computationally intensive because it requires evaluating all possible splits at each step. It often uses methods like k-means or other partitioning techniques to decide the best split. The process can be visualized with a dendrogram, where the root represents the whole dataset and the branches represent the splits. Divisive clustering starts with all data points in a single cluster and splits them iteratively into smaller clusters.

Main steps of the algorithm are:

- **Initialization:** Start with all data points in one cluster.

- **Split Cluster:** Choose a cluster to split and divide it into two sub-clusters.

- **Evaluate Splits:** Evaluate possible splits based on a criterion (e.g., minimizing variance within clusters).

- **Repeat:** Continue splitting clusters until each data point is a singleton cluster or the desired number of clusters is reached.

Main advantages of Hierarchical algorithm are hierarchy and dendogram. Hierarchy basically provides a multi-level hierarchy, which can be useful for understanding data at different levels of granularity. A dendrogram is a tree-like diagram that illustrates the arrangement of clusters formed by hierarchical clustering. Each branch represents a split at various levels of the hierarchy, showing the relationships between clusters. It is used to visualize the process of clustering and the order in which clusters are merged or split.

On the other hand, disadvantages of the divisive hierarchical clustering are computational complexity, sensitivity to noise and no objective function. Computational complexity refers to the amount of computational resources, such as time and space, required to solve a problem as a function of the size of the input. It is typically expressed using Big O notation, which describes the upper bound of the algorithm's growth rate. Understanding computational complexity helps in evaluating and comparing the efficiency of different algorithms. Agglomerative methods can be computationally expensive, especially for large datasets (typically $O(n^3)$).

Sensitivity to noise and outliers refers to how much an algorithm's performance is affected by the presence of erroneous or extreme data points. Algorithms highly sensitive to noise and outliers can produce distorted or incorrect results, while robust algorithms can handle such data more effectively. Techniques like preprocessing, using robust metrics, or employing algorithms specifically designed to mitigate noise and outliers can help manage this sensitivity. Addressing this issue is crucial for ensuring the reliability of data analysis and modeling. Unlike k-means, hierarchical clustering does not have a clear objective function, making the determination of the number of clusters more subjective.

Divisive hierarchical clustering has a range of applications across various fields due to its ability to identify and separate distinct subgroups within a dataset. In market segmentation, it helps businesses understand distinct customer groups, enabling targeted marketing strategies. In bioinformatics, it aids in the classification of genes or proteins with similar functions, contributing to understanding biological processes and disease mechanisms. In document clustering, it can organize large sets of texts into meaningful categories, improving information retrieval and topic modeling. Additionally, divisive clustering is useful in anomaly detection, where it helps isolate outliers from normal data points, enhancing security and fraud detection systems. Its flexibility and interpretability make it a valuable tool in exploratory data analysis and pattern recognition.

**Example Code (Agglomerative Clustering in Python using scikit-learn):**

```
from sklearn.cluster import AgglomerativeClustering
import numpy as np

# Sample data
X = np.array([[1, 2], [1, 4], [1, 0],
        [4, 2], [4, 4], [4, 0]])

# Agglomerative clustering
agg_cluster = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='ward')
agg_cluster.fit(X)

# Cluster labels
labels = agg_cluster.labels_
print(labels)
```

# 7.6 Density Based Algorithms

Density-based clustering algorithms are designed to discover clusters of varying shapes and sizes in data, particularly in situations where clusters may not be well-separated or have different densities. These algorithms work by identifying dense regions of data points, which are separated by regions of lower density. The most commonly used density-based clustering algorithms are DBSCAN (Density-Based Spatial Clustering of Applications with Noise) and OPTICS (Ordering Points To Identify the Clustering Structure).

## 7.6.1 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

### 7.6.1.1 Key Concepts

- **Core Point:** A point is a core point if it has at least a minimum number of neighboring points (MinPts) within a given radius ($\varepsilon$).

- **Border Point:** A point that is not a core point but falls within the $\varepsilon$ radius of a core point.

- **Noise Point:** A point that is neither a core point nor a border point.

### 7.6.1.2 Steps

- **Initialization:** Start with an arbitrary point in the dataset.

- **Check Density:** If the point is a core point, create a new cluster and retrieve all points density-reachable from the core point.

- **Expand Cluster:** Add density-reachable points to the cluster, and repeat the process for all newly added core points.

- **Mark Noise:** Points that are not reachable from any core point are marked as noise.

- **Repeat:** Continue until all points have been visited.

### 7.6.1.3 Parameters

- **$\varepsilon$ (epsilon):** The radius within which to search for neighboring points.

- **MinPts:** The minimum number of points required to form a dense region.

### 7.6.1.4 Advantages and disadvantages

Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is a popular clustering algorithm known for its ability to identify clusters of arbitrary shape and handle noise effectively. Among its advantages, DBSCAN does not require specifying the number of clusters a priori, making it well-suited for exploratory data analysis. It excels in discovering clusters of varying densities and is robust to noise and outliers, as it can classify them as separate noise points. This robustness enhances its application in real-world scenarios where data is often imperfect.

However, DBSCAN also has some disadvantages. The performance of the algorithm depends heavily on the choice of its two parameters: the radius ($\varepsilon$) and the minimum number of points (minPts). Selecting these parameters can be challenging, especially for datasets with varying densities. Additionally, DBSCAN may struggle with high-dimensional data where defining a meaningful neighborhood radius is difficult. It can also be computationally expensive for large datasets due to the need to calculate distances between all data points. Despite these drawbacks, DBSCAN remains a powerful tool for density-based clustering, particularly in applications where noise and irregular cluster shapes are prevalent.

## 7.6.2 OPTICS (Ordering Points To Identify the Clustering Structure)

Ordering Points To Identify the Clustering Structure (OPTICS) is a density-based clustering algorithm designed to overcome some of the limitations of DBSCAN by effectively identifying clusters of varying densities. Unlike DBSCAN, OPTICS does not require a predefined radius ($\varepsilon$) for clusters, allowing it to adapt to local density variations within the data. The algorithm works by ordering points to create a reachability plot, which visually represents the structure of the dataset and highlights areas of different density. This plot helps in understanding the hierarchical relationships between clusters and assists in selecting optimal clustering parameters. OPTICS is particularly valuable in applications where the data has a complex structure with clusters of different shapes and densities, providing deeper insights and more accurate clustering results.

*7.6.2.1 Key Concepts*

- **Core Distance:** The minimum $\varepsilon$ (radius) needed for a point to be a core point.
- **Reachability Distance:** The distance used to maintain the order of points based on density reachability.

*7.6.2.2 Steps*

- **Initialize:** Start with an arbitrary point and compute its core distance.
- **Expand Clusters:** Similar to DBSCAN, expand clusters by retrieving density-reachable points, but store the order and reachability distance.
- **Ordering:** Maintain an order of points based on their reachability distances.
- **Extract Clusters:** Use the ordered list to extract clusters with varying densities.

*7.6.2.3 Advantages and Disadvantages*

Ordering Points To Identify the Clustering Structure (OPTICS) is an advanced density-based clustering algorithm that addresses some of the limitations of DBSCAN. One of its main advantages is that it does not require a fixed radius ($\varepsilon$) for all clusters, allowing it to identify clusters with varying densities more effectively. This adaptability makes OPTICS particularly useful for complex datasets where the density of clusters changes significantly. Additionally, OPTICS produces a reachability plot that provides a visual representation of the clustering structure, helping to identify the optimal clustering parameters and offering insights into the hierarchical relationships among clusters.

However, OPTICS also has its disadvantages. The algorithm can be computationally intensive, especially for large datasets, as it involves sorting points based on their reachability distance, which can be time-consuming. Furthermore, interpreting the reachability plot can be challenging, requiring a good understanding of the data and the algorithm's mechanics to identify meaningful clusters accurately. Despite these challenges, OPTICS remains a valuable tool for density-based clustering, offering greater flexibility and insight into the clustering structure compared to simpler methods like DBSCAN.

*7.6.2.4 Example Implementation (DBSCAN in Python using scikit-learn)*

Here's how you can implement DBSCAN using scikit-learn:

```python
import numpy as np

from sklearn.cluster import DBSCAN

import matplotlib.pyplot as plt


# Sample data
X = np.array([[1, 2], [2, 2], [2, 3],

        [8, 7], [8, 8], [25, 80]])


# DBSCAN clustering
db = DBSCAN(eps=3, min_samples=2).fit(X)

labels = db.labels_


# Plotting the results
plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis')

plt.title('DBSCAN Clustering')

plt.xlabel('Feature 1')

plt.ylabel('Feature 2')

plt.show()
```

Choosing the right parameters (ε and MinPts) is crucial for DBSCAN to work effectively. This often requires domain knowledge or heuristic methods such as:

- **k-distance graph:** Plotting the distance to the k-th nearest neighbor and finding the "elbow" point.

- **Domain expertise:** Understanding the data to make informed choices about the density thresholds.

## 7.6.3 Applications

- **Spatial Data Analysis:** Identifying geographical patterns, such as earthquake epicenters or crime hotspots.

- **Image Processing:** Segmenting images based on color or texture density.

- **Market Analysis:** Finding dense regions of customers in demographic or transactional data.

Consequently, Density-based clustering algorithms like DBSCAN and OPTICS are powerful tools for identifying clusters in data with varying shapes and densities. They are particularly useful when dealing with noisy data or when the number of clusters is not known in advance. By focusing on the density of data points, these algorithms can uncover complex structures that other clustering methods might miss.

# 7.7 Grid Based Algorithms

Grid-based algorithms for clustering partition the data space into a finite number of cells forming a grid structure. These algorithms are particularly efficient for handling large datasets because they reduce the complexity by focusing on the dense regions within the grid, rather than on individual data points. A primary advantage of grid-based clustering is its speed, as it processes data in a fixed number of operations, making it less sensitive to the size of the dataset. Examples include STING (Statistical Information Grid) and CLIQUE, which are used in various fields such as spatial data analysis and data mining.

However, grid-based algorithms also have limitations. Their performance and accuracy depend heavily on the resolution of the grid, which can be challenging to set appropriately. If the grid resolution is too coarse, important patterns may be missed; if too fine, the computation may become unnecessarily complex. Additionally, these algorithms may struggle with high-dimensional data, where defining meaningful grids becomes difficult. Despite these challenges, grid-based clustering remains a powerful tool for specific applications, particularly when dealing with large-scale spatial data.

## 7.7.1 Principles of Grid-Based Algorithms

1. **Grid Representation:** A grid is a spatial representation that divides a space into regular, usually square, cells or tiles. Each cell can represent various types of data or states, depending on the application.

2. **Indexing:** Grid cells are typically indexed using integer coordinates, which makes accessing and updating cells straightforward.

3. **Spatial Partitioning:** Grids partition the space into manageable chunks, enabling efficient spatial queries, collision detection, and neighbor searches.

4. **Efficiency:** Due to their regular structure, grids allow for efficient algorithms with predictable time complexity, often linear or near-linear in terms of the number of cells.

## 7.7.2 Applications of Grid-Based Algorithms

1. **Pathfinding:** Algorithms like *A\** and *Dijkstra's* algorithm use grids to represent the search space for finding the shortest path between points.

2. **Collision Detection:** In physics simulations and games, grids help detect collisions between objects by checking only the relevant cells for potential overlaps.

3. **Image Processing:** Grids are used to represent pixel data in images, enabling efficient image manipulation, filtering, and analysis.

4. **Geographic Information Systems (GIS)**: Grids are used to represent spatial data, such as elevation, land use, and population density, facilitating efficient querying and analysis.

5. **Robotics:** Grids are used in robot navigation and mapping (e.g., occupancy grids in SLAM - Simultaneous Localization and Mapping).

## 7.7.3 Examples of Grid-Based Algorithms

### 7.7.3.1 Pathfinding Algorithms

Pathfinding algorithms are crucial for finding the optimal route between two points in various applications like robotics, gaming, and network routing. They are broadly categorized into uninformed and informed methods. Uninformed algorithms, such as Breadth-First Search (BFS) and Depth-First Search (DFS), explore paths without additional information about the goal, with BFS guaranteeing the shortest path in unweighted graphs and DFS offering a depth-oriented exploration. In contrast, informed algorithms like A\* and Dijkstra's use heuristics or cost functions to enhance efficiency and accuracy. A\* combines the cost to reach a node with an estimated cost to the goal, providing an efficient pathfinding solution, while Dijkstra's algorithm is effective for graphs with known edge weights. The choice of algorithm depends on factors like the complexity of the graph, the need for speed, and the nature of the pathfinding task.

- **A\* Algorithm:** *A\** is a popular pathfinding algorithm that uses a heuristic to guide its search. The grid represents the search space, with each cell corresponding to a possible position. The

algorithm evaluates the cost of moving from the start to the goal through various cells, considering both the actual distance traveled and an estimated distance to the goal.

- **Dijkstra's Algorithm:** Similar to *A\**, but without the heuristic. It explores all possible paths from the start node, ensuring that the shortest path to each cell is found.

### 7.7.3.2 Cellular Automata

- **Game of Life:** This is a cellular automaton where each cell on a grid can be alive or dead. The state of each cell in the next generation is determined by its current state and the states of its eight neighbors.

### 7.7.3.3 Spatial Hashing

- **Collision Detection:** In a physics engine, objects are placed into grid cells based on their positions. To check for collisions, the algorithm only needs to check objects within the same cell or adjacent cells, significantly reducing the number of checks compared to a brute-force approach.

### 7.7.3.4 Image Processing

- **Convolution Operations:** Grids representing images are processed using convolution filters to perform operations like blurring, edge detection, and sharpening. Each pixel value is updated based on a weighted sum of its neighbors.

### 7.7.3.5 Heatmap Generation

- **Data Visualization:** Grids are used to aggregate spatial data into cells to generate heatmaps, which visualize the density or intensity of data points across an area.

## 7.7.4 Grid Types and Variations

Grid types and variations are used in various fields for organizing and analyzing spatial data, optimizing algorithms, and managing resources. Here's an overview of some common grid types and their variations:

1. **Regular Grids**: These grids have cells of uniform size and shape, such as square or rectangular grids. They are commonly used in geographic information systems (GIS) and computer graphics. Regular grids simplify data management and processing but may not adapt well to varying data densities or irregular shapes.

2. **Hexagonal Grids**: These grids use hexagon-shaped cells, providing a more uniform distance between adjacent cells compared to square grids. Hexagonal grids are often used in applications requiring smooth and continuous coverage, such as game maps or spatial analysis.

3. **Adaptive Grids**: Adaptive or dynamic grids change cell sizes based on data density or complexity. These grids have smaller cells in areas with high data density and larger cells where data is sparse. They are useful for applications like climate modeling or adaptive mesh refinement in simulations.

4. **Quadtree Grids**: This hierarchical grid divides the space into four quadrants recursively. Quadtree grids are beneficial for spatial indexing and managing varying levels of detail, often used in computer graphics and spatial databases.

5. **Octree Grids**: Similar to quadtrees, but in three dimensions, octree grids partition space into eight octants. They are used in 3D modeling, volumetric data processing, and spatial indexing for applications such as 3D game engines and simulations.

## 7.7.5 Challenges and Considerations

- **Resolution and Memory Trade-offs:** Higher resolution grids provide more detail but require more memory and processing power.

- **Boundary Conditions:** Handling the edges and corners of the grid can introduce complexities, especially in wrapping or clipping data.

- **Scalability:** For very large grids or real-time applications, algorithms must be optimized for performance, potentially using parallel processing or specialized data structures.

Grid-based algorithms are powerful tools for solving a variety of spatial problems efficiently. Their regular structure allows for straightforward implementation and optimization, making them a go-to choice in fields ranging from computer graphics to robotics. Understanding the principles and applications of grid-based algorithms enables developers and researchers to design effective solutions for complex spatial challenges.

## Check Your Progress

- What is the main difference between agglomerative (bottom-up) and divisive (top-down)?
- Discuss the basic concept of density based algorithm.
- Write any two application areas of grid based algorithms.

## 7.8 Model Based Algorithms

Model-based algorithms for clustering and data analysis involve creating a statistical model that represents the underlying structure of the data. These algorithms assume that the data is generated from a mixture of underlying probabilistic distributions, such as Gaussian distributions in the case of Gaussian Mixture Models (GMMs). The algorithm aims to fit the model to the data, estimating parameters that best describe the clusters or patterns present. Model-based approaches are powerful because they can handle complex data distributions and provide probabilistic cluster assignments, offering a nuanced view of the data. However, they can be computationally intensive and sensitive to initial conditions or assumptions about the data distribution. Despite these challenges, model-based algorithms are widely used in various applications, including pattern recognition, anomaly detection, and complex data analysis, due to their ability to capture and represent intricate data structures.

## 7.8.1 Types of Model-Based Algorithms

### 7.8.1.1 Dynamic Programming (DP)

- **Definition:** DP is a method for solving complex problems by breaking them down into simpler subproblems. It is applicable to problems exhibiting the properties of overlapping subproblems and optimal substructure.

- **Applications:** Resource allocation, route optimization, sequence alignment in bioinformatics.

- **Example:** The Bellman equation in reinforcement learning.

### 7.8.1.2 Kalman Filters

- **Definition:** A recursive algorithm used for estimating the state of a dynamic system from a series of noisy measurements.

- **Applications:** Navigation systems, financial market predictions, control systems.

- **Example:** GPS systems use Kalman filters to predict the current location based on past data.

### 7.8.1.3 Hidden Markov Models (HMMs)

- **Definition:** Statistical models where the system being modeled is assumed to follow a Markov process with hidden states.

- **Applications:** Speech recognition, bioinformatics (gene prediction), financial market analysis.

- **Example:** Viterbi algorithm for decoding the most likely sequence of hidden states.

### 7.8.1.4 Bayesian Networks

- **Definition:** Probabilistic graphical models that represent a set of variables and their conditional dependencies using a directed acyclic graph.

- **Applications:** Medical diagnosis, machine learning, risk analysis.

- **Example:** Diagnosis of diseases based on symptoms.

## 7.8.2 Model-Based Reinforcement Learning (MBRL)

### 7.8.2.1 Model Learning

- **Definition:** In MBRL, the agent learns a model of the environment's dynamics and uses this model to simulate outcomes and make decisions.

- **Techniques:**

  1. Forward Modeling: Predicting the next state given the current state and action.

  2. Inverse Modeling: Determining the action needed to reach a specific state from the current state.

*7.8.2.2 Planning and Control*

- **Definition:** Using the learned model to plan a sequence of actions to achieve a goal.

- **Techniques:**

  1. Model Predictive Control (MPC): Optimizing actions over a finite horizon and applying the first action.
  2. Monte Carlo Tree Search (MCTS): Simulating many possible action sequences and choosing the best one based on expected rewards.

*7.8.2.3 Algorithms*

- **Dyna-Q:** Combines model-free Q-learning with model-based planning by using a learned model to generate synthetic experience.

- **PILCO (Probabilistic Inference for Learning Control):** Uses Gaussian processes to model system dynamics and optimize control policies.

## 7.8.3 Advantages and Challenges

*7.8.3.1 Advantages*

- **Sample Efficiency:** Model-based algorithms often require fewer samples to learn effectively, as they can simulate experiences.

- **Generalization:** A well-learned model can generalize to new situations better than purely model-free approaches.

- **Flexibility:** They can adapt to changes in the environment by updating the model.

*7.8.3.2 Challenges*

- **Model Accuracy:** The effectiveness of model-based algorithms heavily depends on the accuracy of the learned model.

- **Computational Complexity:** Simulating future states and planning can be computationally intensive.

- **Exploration vs. Exploitation:** Balancing the exploration of new strategies and the exploitation of known successful strategies is critical.

## 7.8.4 Applications

- **Robotics:** Path planning and control using dynamic models of the robot and environment.

- **Finance:** Predicting market trends and optimizing trading strategies.

- **Healthcare:** Personalized treatment planning based on models of disease progression.

- **Gaming and Simulations:** AI agents that use models to plan moves and strategies in complex games.

Model-based algorithms are powerful tools for decision-making and prediction across various domains. Their ability to leverage learned or predefined models of systems makes them particularly useful in scenarios where sample efficiency and generalization are crucial. However, the success of these algorithms is contingent upon the quality and accuracy of the models they rely on.

# 7.9 Constraint Based Algorithms

Constraint-based algorithms are designed to cluster data or solve optimization problems by incorporating specific constraints that guide the process. These constraints can be based on domain knowledge, data characteristics, or business rules, allowing the algorithm to focus on feasible and relevant solutions. For example, in constraint-based clustering, constraints such as must-link or cannot-link relationships specify which data points should or should not be clustered together, enabling more meaningful and contextually accurate groupings. These algorithms are widely used in fields like bioinformatics, where constraints may reflect biological relationships, or in customer segmentation, where business rules guide the clustering process. While constraint-based algorithms can produce more relevant results by incorporating external knowledge, they can also be complex to implement and may require careful tuning of constraints to avoid overfitting or missing important patterns.

## 7.9.1 Types of Constraint-Based Problems

### 7.9.1.1 Constraint Satisfaction Problems (CSPs)

- **Definition:** Problems where the goal is to find a set of variable assignments that satisfy all given constraints.

- **Components:**

  1. Variables: The elements that need to be assigned values.
  2. Domains: The possible values that each variable can take.
  3. Constraints: The rules that restrict the values that the variables can simultaneously take.

- **Examples:** Sudoku, n-queens problem, scheduling.

### 7.9.1.2 Constraint Optimization Problems (COPs)

- **Definition:** Problems where the goal is to find an optimal solution from the set of feasible solutions that satisfy all constraints.

- **Components:**

  1. Objective Function: The function to be optimized (minimized or maximized).
  2. Constraints: The conditions that must be met.

- **Examples:** Linear programming, traveling salesman problem, resource allocation.

## 7.9.2 Common Techniques and Algorithms

### 7.9.2.1 Backtracking

- **Definition:** A recursive algorithm for solving CSPs by incrementally building candidates to the solutions and abandoning a candidate as soon as it determines that the candidate cannot possibly be completed to a valid solution.

- **Procedure:**
  1. Choose a variable.
  2. Assign a value to the variable.
  3. Check constraints. If any constraints are violated, backtrack.
  4. Repeat until a solution is found or all possibilities are exhausted.

- **Examples:** Solving puzzles like Sudoku.

### 7.9.2.2 Constraint Propagation

- **Definition:** Techniques used to reduce the search space by deducing variable assignments that must hold for the constraints to be satisfied.

- **Techniques:**
  1. Arc Consistency: Ensuring that for every value of one variable, there is a consistent value in the connected variable.
  2. Node Consistency: Ensuring that every value in a variable's domain satisfies the variable's unary constraints.
  3. Path Consistency: Ensuring that for any pair of variables, every value assignment that satisfies the binary constraint can be extended to a third variable.

- **Examples:** Maintaining arc consistency in scheduling problems.

### 7.9.2.3 Local Search Algorithms

- **Definition:** Techniques that start from an initial candidate solution and iteratively move to a neighboring solution by making small changes, aiming to find a solution that satisfies all constraints or optimizes the objective function.

- **Techniques:**
  1. Hill Climbing: Continuously moving to neighboring solutions with better objective function values.
  2. Simulated Annealing: Similar to hill climbing but allows for occasional moves to worse solutions to escape local optima.
  3. Genetic Algorithms: Using principles of natural selection to evolve solutions over generations.

- **Examples:** Solving large-scale optimization problems like the traveling salesman problem.

*7.9.2.4 Linear and Integer Programming*

- **Definition:** Mathematical methods for optimizing a linear objective function subject to linear equality and inequality constraints.

- **Types:**

  1. Linear Programming (LP): Variables can take any continuous values.
  2. Integer Programming (IP): Variables are constrained to integer values.

- **Techniques:**

  1. Simplex Algorithm: A popular algorithm for solving LP problems.
  2. Branch and Bound: A method for solving IP problems by partitioning the feasible region into smaller subproblems.
  3. Cutting Planes: Techniques for solving IP by iteratively adding constraints to cut off non-integer solutions.

- **Examples:** Resource allocation, network flow optimization.

## 7.9.3 Advantages and Challenges

*7.9.3.1 Advantages*

- **Flexibility:** Can be applied to a wide range of problems with different types of constraints.

- **Exact Solutions:** Capable of finding exact solutions for many problems.

- **Modeling Power:** Can model complex relationships and dependencies between variables.

*7.9.3.2 Challenges*

- **Computational Complexity:** Many constraint-based problems are NP-hard, making them computationally challenging to solve for large instances.

- **Scalability:** Algorithms may struggle with scalability as the number of variables and constraints increases.

- **Handling Non-linear Constraints:** Non-linear constraints can significantly complicate the solution process.

## 7.9.4 Applications

- **Scheduling:** Assigning resources to tasks while satisfying constraints like resource availability, task dependencies, and deadlines.

- **Resource Allocation:** Distributing resources among competing activities while optimizing certain criteria and satisfying constraints.

- **Configuration Problems:** Designing systems or products by selecting and configuring components in a way that satisfies constraints.

- **Artificial Intelligence:** Solving puzzles, games, and logical inference problems.

Constraint-based algorithms are powerful tools for solving a wide variety of problems where constraints define the feasible solutions. By leveraging techniques like backtracking, constraint propagation, local search, and mathematical programming, these algorithms can find solutions that satisfy all constraints or optimize an objective function within a constrained environment. However, the complexity and scalability of these algorithms remain significant challenges, particularly for large-scale and non-linear problems.

## 7.10 Outlier Analysis

Outlier analysis is the process of identifying data points that significantly deviate from the majority of the dataset, which can provide valuable insights into unusual or anomalous behavior. Outliers can arise due to errors in data collection, natural variability, or rare events, and analyzing them helps in understanding these deviations better. Techniques for outlier detection include statistical methods, such as Z-scores and box plots, and machine learning approaches, like isolation forests and one-class SVMs. Effective outlier analysis is crucial in various domains, including fraud detection, quality control, and anomaly detection in systems, as it helps in improving data quality, enhancing model performance, and uncovering hidden patterns or trends. However, distinguishing between genuine outliers and noise can be challenging, requiring careful consideration of the context and domain knowledge.

### 7.10.1 Types of Outliers

*7.10.1.1 Point Outliers*

- **Definition:** Individual data points that are significantly different from the rest of the data.

- **Examples:** A single fraudulent transaction in a dataset of legitimate transactions.

*7.10.1.2 Contextual Outliers*

- **Definition:** Data points that are considered outliers in a specific context but not in others.

- **Examples:** A temperature reading of 25°C might be normal in summer but abnormal in winter for a specific location.

*7.10.1.3 Collective Outliers*

- **Definition:** A group of data points that collectively deviate from the overall data pattern, even if individual points are not outliers.

- **Examples:** A sudden spike in network traffic might indicate a DDoS attack.

## 7.10.2 Techniques for Outlier Detection

### 7.10.2.1 Statistical Methods

- **Definition:** These methods assume a particular statistical distribution of the data and identify outliers based on deviations from this distribution.

- **Techniques:**

  1. Z-Score: Measures how many standard deviations a data point is from the mean. Points with high absolute Z-scores are considered outliers.
  2. Grubbs' Test: Detects outliers in a univariate data set by comparing the deviation of the suspected outlier to the standard deviation of the entire dataset.
  3. Boxplot Method: Uses the interquartile range (IQR) to identify outliers.

### 7.10.2.2 Distance-Based Methods

- **Definition:** Identify outliers based on the distance between data points.

- **Techniques:**

  1. k-Nearest Neighbors (k-NN): Data points with distances significantly larger than their neighbors can be considered outliers.
  2. DBSCAN (Density-Based Spatial Clustering of Applications with Noise): Identifies outliers as points that do not belong to any cluster based on density.

### 7.10.2.3 Clustering-Based Methods

- **Definition:** Use clustering techniques to group similar data points together, and points that do not fit well into any cluster are considered outliers.

- **Techniques:**

  1. K-Means Clustering: Data points far from any cluster centroid can be considered outliers.
  2. Gaussian Mixture Models (GMM): Points with low probability under the model are identified as outliers.

### 7.10.2.4 Model-Based Methods

- **Definition:** Fit a model to the data and identify outliers based on deviations from the model predictions.

- **Techniques:**

  1. Regression Analysis: Points with large residuals (difference between observed and predicted values) are outliers.
  2. Autoencoders: Neural networks trained to reconstruct input data. Points with high reconstruction error are considered outliers.

*7.10.2.5 Machine Learning-Based Methods*

- **Definition:** Use supervised or unsupervised learning algorithms to detect outliers.

- **Techniques:**

  1. Isolation Forest: Constructs an ensemble of trees to isolate data points. Points that require fewer splits to isolate are considered outliers.
  2. Support Vector Machine (SVM): One-class SVM identifies outliers based on the learned boundary of the data distribution.

## 7.10.3 Challenges in Outlier Analysis

*7.10.3.1 High Dimensionality*

- **Problem:** In high-dimensional spaces, the notion of distance becomes less meaningful, making it harder to identify outliers.

- **Solution:** Dimensionality reduction techniques like PCA or t-SNE can be used to project data into lower-dimensional spaces.

*7.10.3.2 Scalability*

- **Problem:** Detecting outliers in large datasets can be computationally intensive.

- **Solution:** Use scalable algorithms like incremental learning or distributed computing frameworks.

*7.10.3.3 Noise*

- **Problem:** Distinguishing between true outliers and noise can be challenging.

- **Solution:** Robust statistical methods and noise-tolerant algorithms can help mitigate this issue.

*7.10.3.4 Interpretability*

- **Problem:** Understanding why a data point is classified as an outlier can be difficult.

- **Solution:** Use explainable models and provide visualization techniques to help interpret the results.

## 7.10.4 Applications of Outlier Analysis

*7.10.4.1 Fraud Detection*

- **Description:** Identifying fraudulent transactions in financial datasets.

- **Techniques:** Statistical methods, clustering-based methods, and machine learning models.

*7.10.4.2 Network Security*

- **Description:** Detecting unusual network traffic patterns that may indicate attacks.

- **Techniques:** Clustering-based methods, distance-based methods, and machine learning models.

*7.10.4.3 Industrial Monitoring*

- **Description:** Identifying anomalies in sensor data from industrial equipment to predict failures.

- **Techniques:** Model-based methods and statistical methods.

*7.10.4.4 Healthcare*

- **Description:** Detecting abnormal patient health indicators to diagnose diseases.

- **Techniques:** Machine learning-based methods and model-based methods.

*7.10.4.5 Environmental Monitoring*

- **Description:** Identifying unusual environmental patterns like sudden changes in temperature or pollution levels.

- **Techniques:** Statistical methods and model-based methods.

Outlier analysis is a vital tool across various domains, helping to identify unusual patterns that could indicate errors, fraud, or novel phenomena. By leveraging a range of statistical, machine learning, and domain-specific techniques, outlier analysis enables the detection and interpretation of these critical data points.

## Check Your Progress

- Give the meaning of constraint-based algorithm.

- Discuss the types of outliers in brief.

- What is the importance of outlier in a machine learning application?

## 7.11 Summary

*Descriptive analytics* involves the process of analyzing historical data to understand what has happened in the past. It employs data aggregation, statistical measures, and data visualization techniques to provide a clear and concise summary of historical performance. By transforming raw data into meaningful insights, it helps identify patterns, trends, and anomalies. Tools like dashboards and reports are commonly used to present the findings. This foundational analysis aids organizations in making informed decisions based on past data.

*Cluster analysis* is a statistical technique used to group similar data points into clusters based on their characteristics. It helps in identifying patterns and structures within the data without predefined labels.

Common methods include K-means, hierarchical clustering, and DBSCAN. These clusters can then be analyzed to uncover insights, such as customer segmentation or anomaly detection. Cluster analysis is widely used in marketing, biology, social network analysis, and many other fields.

*Clustering algorithms* group similar data points into clusters based on their characteristics, facilitating pattern recognition and insights. Common algorithms include K-means, which partitions data into K clusters by minimizing variance within clusters; hierarchical clustering, which builds a tree of clusters using either a bottom-up or top-down approach; and DBSCAN, which identifies clusters based on density and can detect outliers. These techniques are widely used in market segmentation, image analysis, and anomaly detection. Effective clustering helps in understanding the underlying structure of data.

*Partitioning-based* algorithms, such as K-means, divide a dataset into a predefined number of clusters by minimizing the variance within each cluster. These algorithms initialize cluster centroids and iteratively reassign data points to the nearest centroid, recalculating centroids until convergence. They are efficient for large datasets but require specifying the number of clusters in advance. K-medoids, a variation, uses actual data points as centroids for improved robustness against outliers. These methods are widely used in applications like customer segmentation and image compression.

*Hierarchical-based* algorithms build a hierarchy of clusters through either an agglomerative (bottom-up) or divisive (top-down) approach. Agglomerative clustering starts with each data point as its own cluster and iteratively merges the closest pairs until a single cluster remains. Divisive clustering, conversely, starts with one cluster containing all data points and recursively splits it into smaller clusters. These algorithms do not require specifying the number of clusters in advance and produce a dendrogram, a tree-like diagram that illustrates the cluster relationships. Hierarchical clustering is useful for discovering nested patterns in data and is commonly applied in bioinformatics and document clustering.

*Density-based* algorithms, such as DBSCAN (Density-Based Spatial Clustering of Applications with Noise), group data points into clusters based on regions of high density, effectively identifying clusters of arbitrary shape. DBSCAN starts by selecting a random point and expands the cluster by including all points within a specified distance (epsilon) that have a minimum number of neighbors (minPts). Points that do not meet these criteria are labeled as noise or outliers. This method does not require specifying the number of clusters and is effective at handling noise and discovering non-linear cluster structures. Density-based clustering is widely used in spatial data analysis and anomaly detection.

*Grid-based* algorithms partition the data space into a finite number of cells that form a grid structure, and then perform clustering on these cells. One popular method is STING (Statistical Information Grid), which organizes data into hierarchical grid structures and uses statistical information stored in each cell for clustering. These algorithms are efficient in terms of computation and are particularly effective for handling large datasets with multi-dimensional data. Grid-based methods are typically parameter-free concerning the number of clusters, relying instead on grid granularity. They are often used in spatial data mining and geographic information systems.

*Model-based* algorithms assume a specific model for the distribution of data and perform clustering by fitting this model to the data. Gaussian Mixture Models (GMM) are a common example, where clusters are represented as Gaussian distributions and the algorithm estimates parameters such as mean and variance. These methods use statistical approaches to find the best fit and can handle overlapping clusters.

Model-based clustering is effective for probabilistic data analysis and can provide soft assignments of data points to clusters.

*Constraint-based* algorithms incorporate additional constraints or prior knowledge into the clustering process, guiding how clusters are formed. For example, constraints may specify must-link or cannot-link relationships between data points, ensuring that certain points are either in the same cluster or in different clusters. These algorithms can refine clustering results to better align with domain-specific requirements or expert knowledge. Constraint-based clustering is useful in scenarios where relational or contextual information is critical for meaningful cluster formation.

*Outlier analysis* involves identifying and investigating data points that deviate significantly from the rest of the dataset. Techniques include statistical methods like Z-scores and IQR, machine learning approaches such as Isolation Forest and One-Class SVM, and visualization tools like box plots and scatter plots. Outliers may indicate errors, anomalies, or rare but significant events. Analyzing outliers helps in improving data quality and robustness, ensuring accurate insights and predictions. The approach taken often depends on the context and nature of the data.

## 7.12 Terminal Questions

1. What do you understand by descriptive analytics? Discuss the main challenges in descriptive analytics.

2. Explain the difference between K-means and K-medoids clustering algorithms. Write a short note on K-Means clustering algorithm.

3. What are the advantages and disadvantages of using hierarchical clustering compared to K-means clustering?

4. Given a dataset of customer transactions with features like purchase amount, frequency, and category, describe how you would use K-means clustering to segment customers. Include steps for preprocessing, choosing the number of clusters, and interpreting results.

5. Discuss the areas where outliers can be applied efficiently and effcetively.

6. How can outliers impact the performance of a machine learning model, and what are some strategies to mitigate these effects?

7. Consider a dataset of customer ratings for a product. Explain how you would identify and handle outliers that might be the result of fraudulent reviews or user errors. Discuss the choice of methods and the impact on the analysis.

# UNIT-8    PREDICTIVE ANALYTICS – CLASSIFICATION AND PREDICTION

## Structure

## 8.0 Introduction

Predictive analytics is a sophisticated subset of data analytics that focuses on predicting future events using past data. It uses statistical approaches, machine learning algorithms, and data mining processes to evaluate current and historical data, detecting patterns and trends that may be used to make future decisions.  It's a key component of data science and has a wide range of applications across industries, such as finance, healthcare, marketing, and more. Within predictive analytics, **classification** and **prediction** are two fundamental tasks that play a crucial role in making informed decisions.

This Unit Covers several Supervised Learning algorithm used in Classification and Prediction. Using a recursive feature-based data splitting process, a decision tree is a supervised learning technique that builds a model to predict the target variable. Every leaf node represents a value (for prediction tasks) or a class label (for classification tasks) and every interior node represents a feature. Lazy learners—also referred to as instance-based learners—wait to begin learning until a question is posed. Lazy learners

retain the training examples and utilise them straight for prediction, in contrast to eager learners who construct a model during training. The Bayes Theorem is used in Bayesian classification to forecast the likelihood of a class given specific characteristics. It may be applied to both regression and classification tasks and is based on the idea of conditional probability. The model for rule-based categorisation is a collection of if-then rules that link feature values to a target class. The process of examining patterns in the training data yields these rules. One technique for training neural networks, particularly multi-layer perceptrons (MLPs), is backpropagation. The network is made up of several layers of neurones, each of which produces an output by applying a mathematical function to its input. Strong supervised learning methods for regression and classification problems include Support Vector Machines (SVM). SVMs function by maximising the margin between data points by identifying the ideal hyperplane that divides them into distinct classes.

# 8.1 Objective

Objectives of this unit are:

1. To understand the essential ideas and the role of predictive analytics in decision-making.
2. To differentiate categorization and prediction as the two main methods of predictive modeling.
3. To understand several categorization techniques, including Decision Trees, K-Nearest Neighbors, Rule Based Classification, and Support Vector Machines.
4. To learn how to use these algorithms on labeled data to categorize new cases.

# 8.2 Definition

### 8.2.1 Advanced Analytics

This is a data analysis methodology that examines data from many data sources, utilizing statistical techniques, machine learning, deep learning, predictive modeling, and business process automation. There are four types of analytics in Figure 8.1, that come under the umbrella of advanced analytics:

- Descriptive Analytics

- Diagnostic Analytics

- Predictive analytics

- Prescriptive Analytics

The main goal of **Descriptive Analytics** is to provide an overview of previously collected data sets. It also describes what has happened historically. To gain a deeper understanding of the reasons behind occurrences, descriptive analytics analyzes the decisions and outcomes after the fact. Descriptive analytics provides insight into the query of "What happened?".
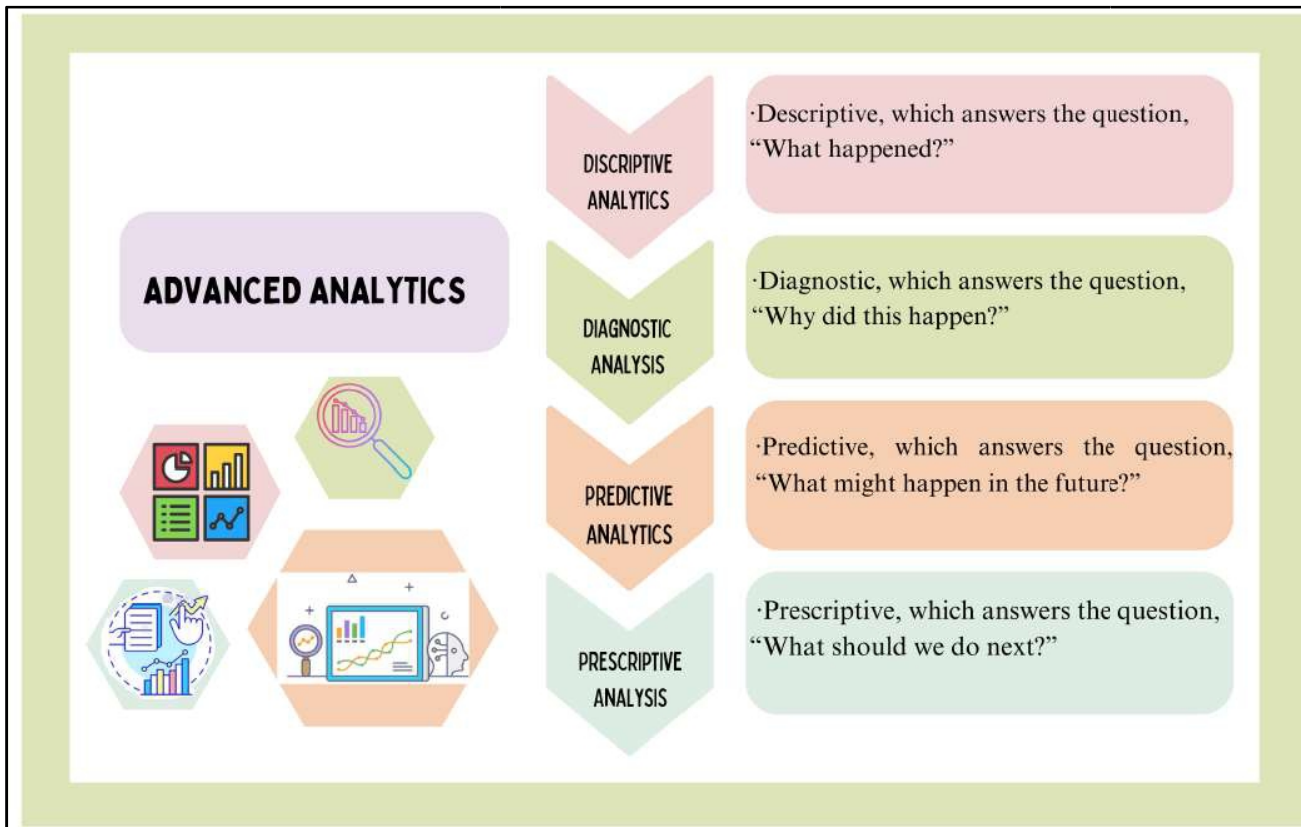


Figure 8.1: Types of Analytics

Causal-effect correlations between variables are identified using **Diagnostic Analytics**. It analyzes the underlying causes of previous occurrences, to help and provide a solution to the question, "Why did it happen?"

Models for **Predictive Analytics** use statistical techniques to project possible outcomes. It provides a solution to the query "What is likely to happen?" by predicting future patterns based on previous data. **Prescriptive Analytics** offers guidance on how to proceed to get desired results. By offering particular recommendations based on the predictions generated by the predictive models, it assists in providing a response to the question, "What should we do?"

**Predictive Analytics:** Predictive analytics determines the probability of future events based on past data by utilizing statistical algorithms, machine learning methods, and historical data. In other words,

the process of extracting data from recent and historical datasets using modeling and statistical approaches to predict probable future events and trends is known as predictive analytics. Predictive analytics leverages prediction techniques and classification to derive insights and forecast future trends, aiding decision-making across various industries.

As a technique in predictive analytics, **predictive modeling** is the process of creating models that look for underlying patterns in both historical and present datasets and estimate the likelihood of a given result using data mining and statistics. Data collection is the first step in the predictive modeling process. After that, a statistical model is created, predictions are produced, and the model is updated as new data becomes available.

**Prediction** in predictive analytics refers to the process of forecasting the outcome or value based on historical data

**Classification:** The process of assigning labels or categorizing each instance, record, or data item in a dataset according to its features or attributes is known as classification in data mining. Predictive modeling techniques such as classification use categories/ classes as the output variable, such as "spam" or "not spam," "diseased" or "not diseased," and so on. Detecting the correct class to which new/ unseen observations belong is the primary objective of classification.  The two types of classification techniques are multi-class and binary classification as shown in Figure 8.2.
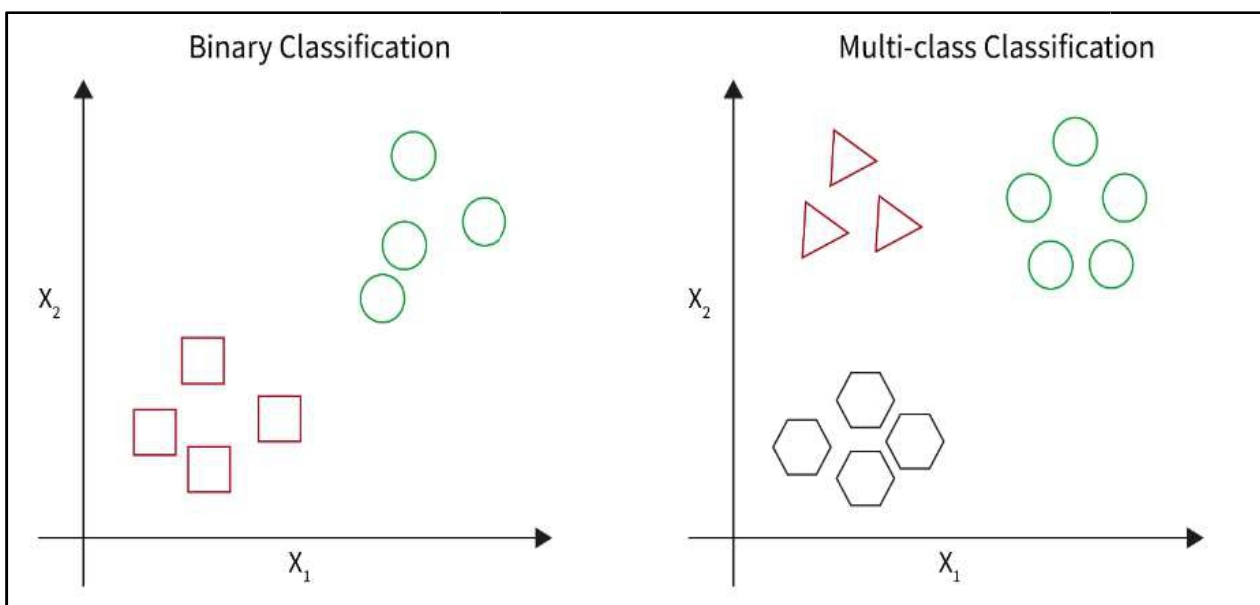


Figure 8.2: Types of Classification

[Source: https://www.geeksforgeeks.org/one-vs-rest-strategy-for-multi-class-classification/]

133

Labels like "fraudulent" or "non-fraudulent" are used in cases of binary categorization. Labeling more than two classes—such as cheerful, neutral, or sad is referred to as multi-class categorization.

**8.2.2 Applications of Predictive Analytics:**

The major applications of predictive analysis are given as follows.

1. Finance: Credit scoring, fraud detection, stock market prediction.

2. Healthcare: Disease prediction, patient risk assessment.

3. Marketing: Customer segmentation, churn prediction, targeted advertising.

4. Manufacturing: Predictive maintenance, quality control.

5. Retail: Inventory management, sales forecasting.

6. Sports: Performance analysis, game outcome prediction.

**8.2.3 Data Mining Classification Categorization:**

Based on their methodology, level of complexity, and level of performance, classification algorithms come in several varieties. There are some of the basic data mining classification categories:

- **Decision Tree-Based Classification**: This kind of algorithm creates a model of decisions and their potential outcomes in the form of a tree. Decision trees are a common solution for categorization issues because they are simple to understand and analyze.

- **Bayesian Classification:** This technique predicts the probability of each class label based on the observable characteristics by using Bayes' theorem. Using Bayesian classification on unclear or partial data is quite helpful.

- **Rule-Based Classification:** An observation's class label is determined by the classification algorithm based on a set of rules. Usually, the rules are written as IF-THEN statements, where each phrase denotes a condition and an appropriate response.

- **Neural Network-Based Classification:** This approach creates a mapping between the input features and the output class labels by using a network of connected nodes or neurons. Complex and nonlinear interactions between features and class labels are well-suited for neural networks.

- **SVM:** For binary and multi-class classification tasks, one of the best supervised learning algorithm is called SVM. The primary concept of Support Vector Machines (SVM) is determining the appropriate hyperplane in a high-dimensional space to divide data points of various kinds.

- **Instance-Based Classification:** Classification algorithms that categorize fresh, unknown examples based on a collection of training instances are known as instance-based classification algorithms. The similarity between the characteristics of the training cases and the new instances serves as the basis for categorization.

When fresh data is provided, instance-based classification determines the label of the closest data point based on similarity (or distance) from the training data that has been stored in memory. This approach mainly uses similarity or distance to predict the labels of fresh data. It is also known as memory-based learning. It is a type of learning algorithms that compares newly observed issue cases with training examples that have been kept in memory, rather than carrying out explicit generalization.

**Advantages and Applications:** The benefits of instance-based learning include the following:
**(i) Flexibility:** Since it isn't dependent on a pre-existing model, it can quickly adjust to changes.
**(ii) Ease of Implementation:** The method is simple to use and understand.
**(iii) No Training Step**: Since every instance represents itself, there is no need for a training step.

**Applications:** Computer Vision, Image Identification, And Recommendation Systems, Among Other Fields.

**Challenges and Limitations:** Instance-based learning has some drawbacks also, given as follows:

**(i)** **Performance:** The quality of the dataset has a significant impact on performance.
**(ii)** **Time and Space Complexity:** The algorithm, particularly for larger datasets, can be computationally expensive and require substantial storage space.
**(iii)** **Sensitivity to Irrelevant Features:** It is affected by irrelevant features, which may cause misclassification.

- **Ensemble-Based Classification:** This method enhances the overall accuracy and resilience of the classification model by combining the predictions of several classifiers. Methods for groups include stacking, boosting, and bagging.

In data mining, an ensemble classifier is a technique that integrates many machine learning models to increase prediction accuracy and overall performance. The key concept of ensemble techniques is that the group may outperform any single model by leveraging the capabilities of several different models.

**Types of Ensemble Techniques:** There are a few main categories of ensemble techniques, each having a unique way of mixing models.

**Bagging Method:** Bagging is the process of training various model iterations on various subsets of the training set. By using replacement sampling to sample the original dataset at random, these subsets are produced (bootstrap sampling). The final output is generated by averaging (for regression) or voting on the predictions made by each individual model after it has been trained individually.

**Boosting:** Boosting is the process of training models one after the other, with the goal of fixing the mistakes caused by the earlier models. In boosting, models are constructed iteratively with the goal of enhancing performance on challenging scenarios, as opposed to being trained singly

**Stacking:** The stacking (also known as stacked generalization) method entails training a number of models (also known as base learners) and then using a different model (also known as a meta-learner or meta-classifier) to combine the predictions of the base learners. The final prediction is made by the meta-learner using the predictions of the basic models, which were trained on the original dataset.

**Voting**: Using a majority vote for classification or an average vote for regression, voting ensembles aggregate the predictions of several models. Voting comes in two primary forms

**(i)**     **Hard Voting** Each model casts a vote for a class label; the final forecast belongs to the class that receives the majority of votes.

**(ii)**     **Soft Voting:** Every model forecasts a class's probability; the class with the greatest average probability across all models is selected.

**Advantages and Disadvantages:**

**Advantages**

**(i)**     **Accuracy**: When compared to individual models, it frequently offers superior prediction performance.

**(ii)**     **Flexibility:** Able to mix many model kinds, even ones with disparate advantages and disadvantages.

**(iii)**     **Robustness:** Less susceptible to overfitting, particularly in the context of bagging techniques.

**Limitations:**

**(i) Complexity:** Ensemble approaches might be more difficult to comprehend and resource-intensive because to their increased complexity and computing demands.

**(ii) Longer Training Times:** It might take a while to train several models, particularly when using techniques like stacking or boosting.

**(iii) Overfitting Risk in Specific Situations:** Certain ensemble techniques, such as boosting, might nevertheless overfit the training set if they are not properly maintained.

---

*Check Your Progress 1.*

Q1. Differentiate descriptive, predictive, and prescriptive analytics.

Q2. Write some examples of descriptive, predictive, and prescriptive analytics.

Q3. What are the challenges of instance-based learning?

Q4. Differentiate hard booting and soft booting.

Q5. What are the advantages of ensemble learning?

---

# 8.3 Decision Tree-Based Classification

Decision trees are supervised learning techniques for regression and classification in data mining. There is a tree that helps in decision-making. The decision tree is often used to create models for classification or regression that resemble trees. The data collection is divided into more manageable groupings to create the decision tree. This tree is a collection of nodes and branches. The tree starts with the root node which is the first node and it may be divided into another number of nodes at the end of the Decision Tree contains only leaf nodes which do not further divide into any child node.

Root nodes might be further divided into some other nodes. These nodes and root nodes are called decision nodes because they represent some condition and on the basis of condition, we take any decision. There are some outcomes that are represented by the child node. so the decision node will have two splits at the absolute least. The leaf nodes stand for categorization or a decision. The best predictor, often referred to as the root node, is connected to the leaf nodes, which are the highest decision nodes in a tree. Decision trees are adaptable enough to manage data that is both categorized and numerical.

Leaf nodes are represented by ovals, and internal nodes are represented by rectangles, as shown in Figure 8.3. Decision tree methods can be used for binary classification as well as multi-class classification in which every internal node branches into two or more additional nodes.
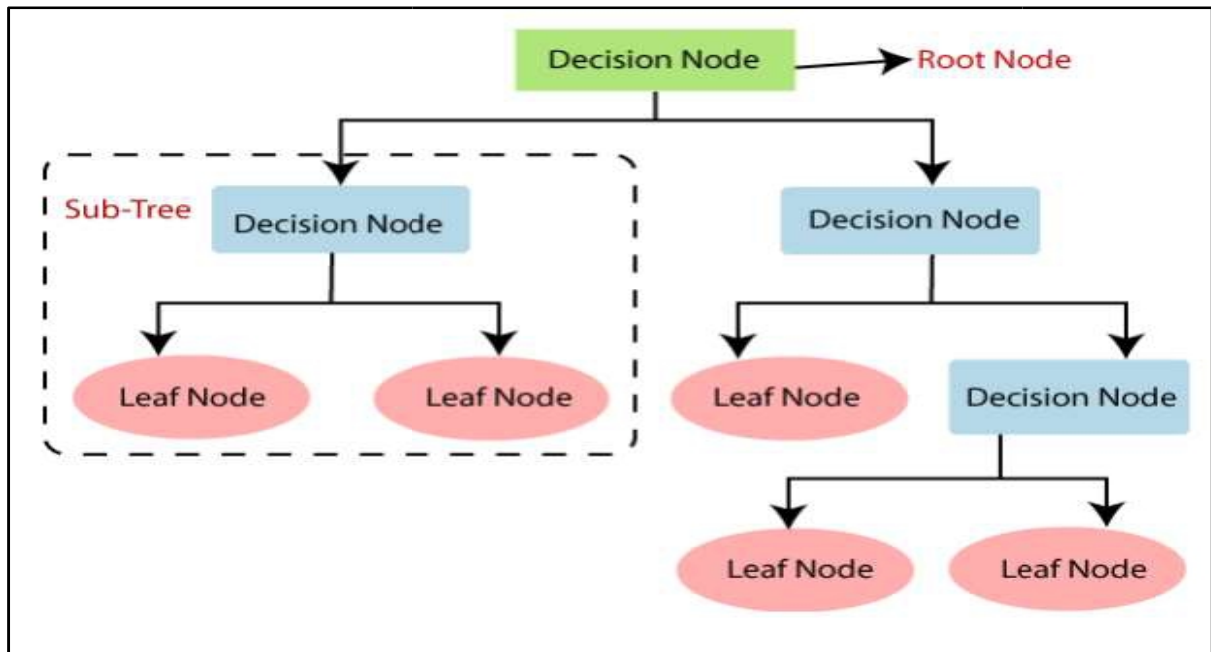


Figure 8.3: Decision Tree

[Source: https://www.javatpoint.com/machine-learning-decision-tree-classification-algorithm]

Now, the next point comes to mind, "How are decision trees used for categorizing objects?" To ascertain if a given tuple (X) is a member of a certain class, its attribute values are compared to those in the decision tree. After leaving the root, the data moves through several intermediary nodes to a leaf node, where the class prediction for that tuple is stored. Decision trees may be used to create classification rules with minimal effort.

### 8.3.1 Steps of ID3 Decision Tree Algorithm:

The initial set S serves as the root node at the start.

- Entropy (H) and information gain (IG) are calculated for each iteration of the algorithm that goes through a highly neglected characteristic in the collection S.

- The property with the highest information gain or the lowest entropy is then chosen.

- After that, a subset of the data is obtained by dividing the set S by the chosen attribute.

- The process keeps repeating on every subset, considering only traits that have never been chosen before.

## 8.3.2 Measures of Attribute Selection

Selecting which attribute to put at the root or at various levels of the tree as internal nodes is a difficult step if the dataset has N attributes. The problem cannot be resolved by choosing any node at random to be the root. A haphazard approach might lead to poor and inaccurate outcomes. The first problem that emerges while creating a decision tree is figuring out which attribute is ideal for the root node and its child nodes. In order to address these issues, a method known as attribute selection measure, or ASM, has been developed. We can quickly choose the ideal attribute for the tree's nodes using this measurement. For ASM, there are two widely used methods, which are:

**1. Information Gain:** Following the segmentation of a dataset according to a feature, information gain is measured as changes in entropy. It determines the amount of knowledge a feature gives us about a class. We divide the node and create the decision tree based on the information gain value. An algorithm using a decision tree is always

| Information Gain= Entropy(S)- [(Weighted Avg) *Entropy(each feature)      eqs(1) |
|---|

**Entropy:** Entropy is a metric to measure the impurity in a given attribute. It specifies randomness in data. Entropy can be calculated as:

| Entropy(s)= -P(yes)log2 P(yes)- P(no) log2 P(no)      eqs(2) |
|---|

Where *pi* is the probability of an instance being classified into a particular class.

- S= Total number of samples
- P(yes)= probability of yes
- P(no)= probability of no

2. **Gini Index:**

- Gini index is a measure of impurity or purity used while creating a decision tree in the CART(Classification and Regression Tree) algorithm.

- An attribute with the low Gini index should be preferred as compared to the high Gini index.

- It only creates binary splits, and the CART algorithm uses the Gini index to create binary splits.

- Gini index can be calculated using the below formula:

$$\text{Gini Index} = 1 - \sum_j P_j^2 \qquad\qquad\qquad eqs(3)$$

Where, $pi$ is the probability of an instance being classified into a particular class.

Since decision tree classifier construction does not need any prior knowledge of the domain, it Is well-suited to discovery through exploration. Decision trees do not have issues with high-dimensional data. Most people easily understand the natural tree structure they apply to display learned knowledge. The learning and categorization procedures of decision tree induction are simple and rapid. Decision trees are often quite accurate in classifying data. However, the result might be impacted by the quality of the data that is accessible. Molecular biology, industry, economics, and medicine are among the sectors that have used decision tree induction methods for classification.

**8.3.3 Advantages of Decision Trees:**

- It is not essential to scale information while utilizing a decision tree.

- The process of creating decision trees is not hampered by data that include missing values either.

- A decision tree model can do the legwork for you; it is automatic and easy to explain to the technical team and stakeholders.

- Decision trees need less pre-processing effort than alternative techniques.

- It is not essential to normalize the data being analyzed in order to use a decision tree.

---

*Check Your Progress 2.*

Q1. Define Entropy, Information Gain, and Gini Index.

Q2. What are the advantages of decision tree?

Q3. Differentiate decision node and leaf node in a decision tree.

---

# 8.4 Lazy Learners

One kind of machine learning algorithm called a lazy learner waits to generalize from the training set of data until the user asks a question of the system. When a prediction is requested, lazy learners store the

training data and execute computation, in contrast to eager learners, who build a model based on the training data before receiving any queries.

### 8.4.1 Characteristics of Lazy Learners:

**1. No Immediate Training-** In the training phase, lazy learners do not construct a model. They employ the training data straight away to create predictions after storing it.

**2. High Query**-Time Complexity: Lazy learners can be computationally costly at prediction time since processing is postponed until a query is made.

**3. Easy to Implement**- As these algorithms need less computation during the training phase, they are simple to implement.

**4. Memory-Intensive-** Keeping track of every training set can take a lot of memory, particularly when dealing with big datasets.

### 8.4.2 Bayesian Classification:

The foundation of Bayesian classification is the Bayes Theorem, which expresses an event's probability in terms of previously known circumstances that may be connected to it. Statistical classifiers that forecast the likelihood that an instance falls into a specific class are known as Bayesian classifiers.

**Bayes' Theorem:**

$$P(A|B)=P(A\cap B)P(B)=P(A)\cdot P(B|A)P(B)$$

where:

$P(A)$= The probability of A occurring

$P(B)$= The probability of B occurring

$P(A|B)$=The probability of A given B

$P(B|A)$= The probability of B given A

$P(A\cap B))$= The probability of both A and B occurring

### 8.4.3 Types of Bayesian Classifiers:

**1. Naive Bayes:** Based on the class label, this model assumes that the characteristics are independent. Although the strong independence assumption is rarely valid in real-world data, it makes the calculation simpler          and          performs          well          in          many          actual          cases.

**2. Gaussian Naive Bayes:** This model assumes a normal (Gaussian) distribution for the continuous features.

**3. Multinomial Naive Bayes:** Usually applied to discrete data, where features denote word frequency, such in text classification.

**4. Bernoulli Naive Bayes:** Used for binary/ boolean features.

### 8.4.4 Benefits of Bayesian Classification:

1. Easy to Use and Quick: Naive Bayes classifiers are quick to learn and need little processing power. 2. Functions Well with Small Data: Surprisingly good performance may be achieved even with a tiny quantity of data.

3. Handles Irrelevant characteristics: Irrelevant characteristics have no appreciable impact on Naive Bayes' performance.

### 8.4.5 Drawbacks of Bayesian classification:

1. Strong Independence Assumption: If the feature independence assumption is not satisfied, Naive Bayes performance may suffer.

2. Zero Probability Problem: The probability estimate for a class and feature value will be zero if they never occur together in the training data. Laplace smoothing is one approach that may be used to reduce this.

### 8.4.6 Applications of Bayesian Classification:

1.  Spam filtering: Recognizing spam emails by looking for certain terms in them.

2.  Text Classification: Sorting papers or news items into different categories.

3.  Medical diagnosis:  Making a knowledgeable assumption about the possibility of a disease based on test results and patient symptoms.

4.  Recommendation systems: Using past data to forecast user preferences.

### 8.4.7 Combining Lazy Learning with Bayesian Classification:

While lazy learners and Bayesian classifiers are conceptually different, they can sometimes be combined. For example, in a k-NN classifier, Bayesian inference could be used to weigh the neighbors differently based on their likelihoods, leading to a more informed prediction.

In summary, lazy learners, like k-NN, and Bayesian classifiers, such as Naive Bayes, offer unique approaches to prediction and classification tasks. Lazy learners are straightforward and computationally

intensive at query time, whereas Bayesian classifiers provide a probabilistic framework that is simple and often effective, especially in scenarios where the independence assumption holds.

---

**Check Your Progress 3.**

Q1. What types of problems are best suited for lazy learning algorithms?

Q2. Write different types of Bayesian classifiers.

Q3. What are advantages of Bayesian classifiers?

Q4. Write the limitations of Bayesian classifiers.

---

# 8.5 Rule-Based Classification

In data mining, rule-based categorization refers to a process where class decisions are made using different "if...then... else" rules. As a result, we describe it as a classification type controlled by an IF-THEN rule set. An IF-THEN rule is expressed as follows: "IF condition THEN conclusion."

**8.5.1 IF-THEN Rule:** We may divide the definition of the IF-THEN rule into two sections:
• **Rule Antecedent:** This is where the rule's "if condition" comes in. The position of this part is on the left-hand side, or LHS. The antecedent can contain one or more features as conditions, using the logic AND operator.

• **Rule Consequent:** This is located on the right-hand side of the rule. The class prediction makes up the rule consequent.

**Example:**

R1: (Give Birth = no) $\wedge$ (Can Fly = yes) $\rightarrow$ Birds

R2: (Give Birth = no) $\wedge$ (Live in Water = yes) $\rightarrow$ Fishes

R3: (Give Birth = yes) $\wedge$ (Blood Type = warm) $\rightarrow$ Mammals

R4: (Give Birth = no) $\wedge$ (Can Fly = no) $\rightarrow$ Reptiles

R5: (Live in Water = sometimes) $\rightarrow$ Amphibians

**8.5.2 Evaluation of the Rule:**

Two criteria are used in rule-based categorization in data mining to get the rules. These are the following:
• **Rule Coverage:** The portion of data that meets a given rule's prerequisite requirements is referred to as

the rule's coverage. We can evaluate this by dividing the number of records that fulfill the rule (n1) by the total number of records (n)

Coverage(R) = n1/n.

• **Rule accuracy:** The percentage of records that meet a rule's consequent values and antecedent requirements is known as the rule's accuracy. This may be computed by dividing the number of records fulfilling the consequent values (n2) by the number of records satisfying the rule (n1). Accuracy (R) = n2/n1.

Usually, we multiply them by 100 to get their percentage. We take this action to simplify these terms and their meanings for everyone.

### 8.5.3 Properties of Rule-Based Classifiers

In data mining, rule-based categorization has two important characteristics. These are:

• Rules might not be exclusive of one another;

• Rules might not be comprehensive.

Rules might not always be incompatible with one another. Since the dataset generates a large number of rules, several of them may fulfill the same data record. This means that the rules are not exclusive of one another. We are not able to select classes that cover distinct data segments on separate rules since the rules do not conflict with one another. However, this was our primary goal. Thus, we have two options for resolving this issue:

Using **an ordered collection of rules** is the first technique. We establish priority orders by orting the rules. This set of ordered rules is thus referred to as a decision list. Hence, the last class is the one with the greatest priority rule. **Assigning votes** to each class based on their weights might be the second option. Thus, the collection of rules is still unordered.

### 8.5.4 Implementation of Rule-Based Classifier

There is two different kinds of methods for implementing rule-based classification in data mining. They are:

- Direct Method

- Indirect Method

### 8.5.4.1 Direct Method

The direct method for rule-based classification in data mining contains algorithms that extract rules directly from the dataset.

The prominent among these are:

- 1R Algorithm

- Sequential Covering Algorithm

### 1R Algorithm

The Learn One Rule Algorithm is alternative name for the 1R Algorithm. It is predicated as an easy categorization principle. As a result, it's among the easiest algorithms. We create rules in this algorithm to test every characteristic. When an attribute contains a range of values, we use the frequency of occurrences of each class value to create rules. By doing this, we can determine its frequency and identify the class with the highest frequency. Next, we designate that class for the pair of attribute values. Next, we determine each attribute's mistake rate and choose the one with the lowest error rate. The 1R approach is not without flaws; noise sensitivity and overfitting are major risks. Thus, in order to solve these issues, we employ the Sequential Covering Algorithm.

### Sequential Covering Algorithm

The sequential covering algorithm is the most widely used algorithm in rule-based classification in data mining. This algorithm defines rules to cover the maximum data records of one class and no data records of other classes. In these algorithms, rules are grown sequentially, one at a time. After a rule is grown, all data records covered by that rule are eliminated, and then this process keeps on repeating for the rest. This process will stop when we reach the **stopping criteria**,

### 8.5.4.2 Indirect Method

We construct rules from many other classification models in the indirect technique of rule-based classification in data mining. Among the well-known categories from which we get the rules are neural networks and decision trees.

### 8.5.5 Benefits of Using a Rule-Based Classifier

There are several advantages associated with using a rule-based classifier:

**1) Transparency** - Users, analysts, and data scientists participating in decision-making processes may readily learn the reasoning behind specific classifications produced by rule-based classifiers since each decision they make adheres to pre-established logic.

**2) Accuracy -** Due to its capacity to understand complex correlations between the variables and features found in those datasets, rule-based classifiers have a tendency to achieve high accuracy levels when appropriately trained using pertinent datasets.

**3) Flexibility -** Unlike other techniques, such as neural networks, where even little adjustments may necessitate retraining whole models from scratch, rules may be readily changed as needed without requiring substantial modifications to the underlying algorithms.

**4) Interpretability-** While these systems base every judgment on pre-established reasoning, analysts and data scientists who participate in the decision-making process will find it simpler to understand the reasoning behind the classifications that were generated.

**5) Scalability-** Rule-based classifiers are easily scalable to handle enormous datasets. This is due to the fact that, regardless of the amount or complexity of the dataset, they work based on a predetermined set of criteria.

**6) Explainability -** In fields like banking and healthcare where transparency is essential, a rule-based classifier's explainability makes it an excellent tool for regulatory compliance.

**7) Speed-** Because rule-based classifiers depend on pre-established rules rather than complex mathematical models, they are often faster than other machine learning methods. This makes them perfect for real-time applications.

8.5.6 Application of Rule Based Classification

**Credit Scoring:** Rule-based classifiers can be used to assess creditworthiness of individuals or businesses by analyzing factors such as income, credit history, and debt-to-income ratio.

**Predictive Maintenance:** By analyzing equipment data, rule-based systems can predict when maintenance is needed before a breakdown occurs, reducing downtime and improving efficiency

**Spam Filtering:** Rule-based classifiers are commonly used in email filters to detect and block unwanted spam messages based on certain criteria such as keywords or sender information.

**Quality Control:** In manufacturing settings, rule-based systems can analyze product quality data to identify defects and improve production processes.

**Healthcare Industry**- Rule-based classifiers have been successfully applied for disease diagnosis

**Finance Industry-** Rule-based classifiers have been widely adopted for fraud detection.

---

*Check Your Progress 4*

Q1. What are the key components of a rule?

Q2. What is rule coverage?

Q3. What are different methods for implementing rule-based classifiers?

Q4. Write major applications of rule-based classifiers.

---

# 8.6 Classification By Back-Propagation

Back-propagation, is a popular neural network training approach, to classification tasks. This methodology is based on the training of feedforward neural networks, particularly multi-layer perceptrons (MLPs), to classify data into distinct categories or classes. In layman's term, Classification using back-propagation involves training a neural network to categorize input data into specific classes by adjusting the network's weights through the back-propagation algorithm. In this process, the neural network first takes in input data and passes it through multiple layers, each comprising neurons that process the data using weighted connections. The output layer then produces a prediction, typically in the form of probabilities for each class. The accuracy of this prediction is assessed by comparing it to the actual class labels using a loss function, such as cross-entropy loss. Back-propagation then calculates the gradients of the loss function concerning each weight in the network by propagating the error backward from the output layer to the input layer. This gradient information is used to update the weights in a direction that minimizes the loss, thereby improving the network's accuracy over time. This process is repeated across many iterations, or epochs, allowing the network to learn complex patterns in the data and make increasingly accurate classifications. Back-propagation is foundational in training neural networks, particularly in deep learning, and is widely used in applications ranging from image recognition to speech processing.

## 8.6.1 Back-propagation

Neural network training algorithms employ the supervised learning technique known as back-propagation. To update the weights and reduce prediction error, it propagates the error from the output layer back through the network layers. Backpropagation gets its name because these changes are propagated "backward," from the output layer via each hidden layer and back to the first hidden layer.

147

The learning process will eventually come to an end when the weights converge, although this is not always assured.

### 8.6.2 The working of Back-propagation:

*Feedforward Phase:* A multilayer feed-forward neural network is used to learn using the backpropagation technique. In this, the input data moves layer by layer through the network. Every layer is made up of neurons, or nodes, except the input layer, every node applies a non-linear activation function and the weighted sum of the inputs to generate an output.



Figure 8.4: Neural Network

[Source: https://www.geeksforgeeks.org/artificial-neural-networks-and-its-applications/]

The inputs match the characteristics that are assessed for every training sample. The neurons in the **input layer** receive the inputs concurrently. Additionally, a second layer of neuron-like units called a **hidden layer** receives the weighted outputs of these units concurrently. An additional hidden layer can get its weighted outputs from the hidden layer, and so on. The **output layer** units, which provide the network's prediction for specific samples, get the weighted outputs of the final hidden layer as input.

The **activation function** decides whether a neuron should be activated or not by calculating the weighted sum and further adding bias to it. The purpose of the activation function is to introduce non-linearity into the output of a neuron. There are several activation functions are used like linear, ELU,

148

ReLU, sigmoid, softmax, tanh, etc. but the softmax function is commonly used by the output layer in classification tasks to transform the raw outputs into probabilities for each class.

**Loss Function:**

A loss function (such as cross-entropy loss for classification) is used to compare the network's output to the actual labels, or the predicted output. The weights are adjusted for each training sample to minimize the mean squared error between the network's prediction and the actual class. The discrepancy between the expected and actual outputs is quantified by the loss function. A set of training samples is used in the backpropagation process repeatedly to learn by comparing each sample's predicted class label with the actual known class label.

Binary Cross Entropy is calculated as:

$$BCE = -\left( y * \log(y\_pred) + (1 - y) * \log(1 - y\_pred) \right) \qquad \text{eqs(4)}$$

Where:

BCE: Binary Cross Entropy

y: True label (either 0 or 1)

y_pred: Predicted probability (between 0 and 1)

log: Natural logarithm (usually base-e logarithm)

**Backpropagation of Errors:**

Errors are propagated from the output layer back to the input layer via the network in back- propagation. Using the calculus chain rule, it determines the gradient of the loss function about each weight in the network. Layer by layer, working backward from the output layer, the error is computed. How much each weight contributes to the mistake is found by calculating the gradients of the loss function concerning the weights.

**Revising the Weights:**

After the gradients are computed, an optimization procedure, usually gradient descent, is used to update the weights. A learning rate is frequently used to regulate the magnitude of the weight updates as the weights are changed in a direction that lowers the loss function.

**The Iterative Method:**

Until the loss converges to a minimum or fails to improve, the feedforward and back-propagation stages are repeated several times, or epochs. The neural network learns the ideal weights to reduce the loss function and raise classification accuracy during this process.

### 8.6.3 Applications of Back-propagation Classification

Neural networks frequently employ back-propagation for various kinds of classification problems, such as:

**Images Classification:**

In image classification applications, neural networks that use back-propagation are frequently used to train the network to identify objects or patterns in pictures. For tasks like as categorising handwritten digits or detecting objects in photos, Convolutional Neural Networks (CNNs) are frequently trained using back-propagation. CNNs consist of numerous layers of convolution and pooling.

**Text Categorisation:**

Text data classification tasks including text classification, sentiment analysis, and spam detection are performed using neural networks. Back-propagation through time (BPTT), a variation of the back-propagation method, is used to train recurrent neural networks (RNNs) and long short-term memory (LSTM) networks, which are made to handle sequential input.

**Speech Recognition:**

Neural networks are trained to classify audio input into phonemes or words in speech recognition. Back-propagation helps in fine-tuning the weights to increase recognition accuracy in networks such as CNNs and RNNs.

**Medical Diagnosis:**

Additionally, neural networks are used in medical diagnostics, where they are trained to identify diseases or conditions based on genetic information, medical images, or patient data.

---

*Check Your Progress 5.*

Q1. What are different types of layers in a neural network?

Q2. What is the role of activation function in neural networks?

Q3. Define backpropagation error.

Q4. What are the major applications of backpropagation method?

---

# 8.7 Support Vector Machine

A supervised learning machine learning approach called "Support Vector Machine" (SVM) can be used for problems including either regression or classification. But the majority of its applications are in classification tasks, including text categorization. Each data point is plotted as a point in n-dimensional space (where n is the number of features) using the SVM method. The value of each feature is represented by a specific coordinate. Next, we perform classification by identifying the ideal hyper-plane that effectively distinguishes between the two classes.

**Hyperplane: A decision boundary that divides several classes of data points in an SVM is called a hyperplane. As an example, the hyperplane is a line in two dimensions and a plane in three dimensions. The main objective of the SVM is to find out the ideal hyperplane to maximize the margin between the classes. The gap between the closest data points from either class and the hyperplane is known as the margin.**

It increases the difference between the nearest data points of different classes. If the hyperplane is a line in a two-dimensional space or a plane in an n-dimensional space, it is determined by the number of features in the input data. The algorithm finds the optimal decision boundary between classes by maximizing the margin between points since numerous hyperplanes might be identified to differentiate classes. Support vectors are the lines that run parallel to the ideal hyperplane and pass across the data points that establish the maximum margin.

• The three hyperplanes (A, B, and C) are shown in Fig. 8.5(a) -8.5(c). In Fig. 8.5(a) Hyper-plane "B" performed an excellent job in this case.
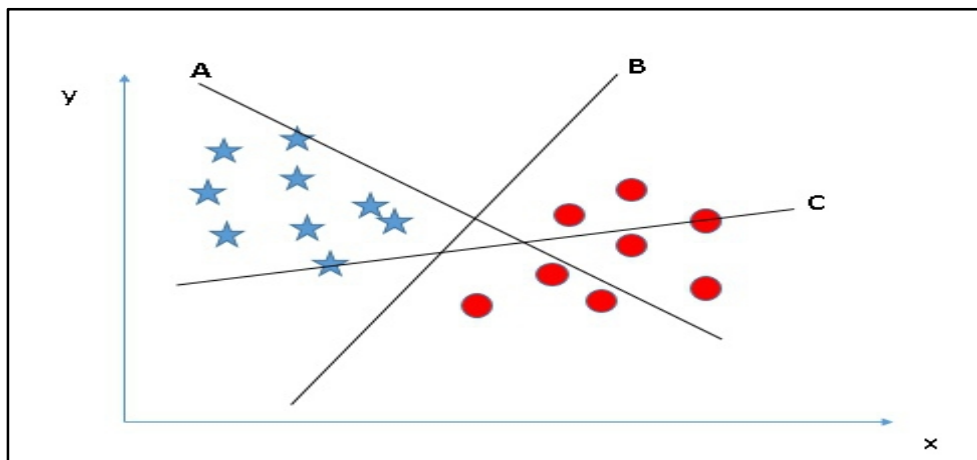


Figure 8.5(a) : Hyperplane

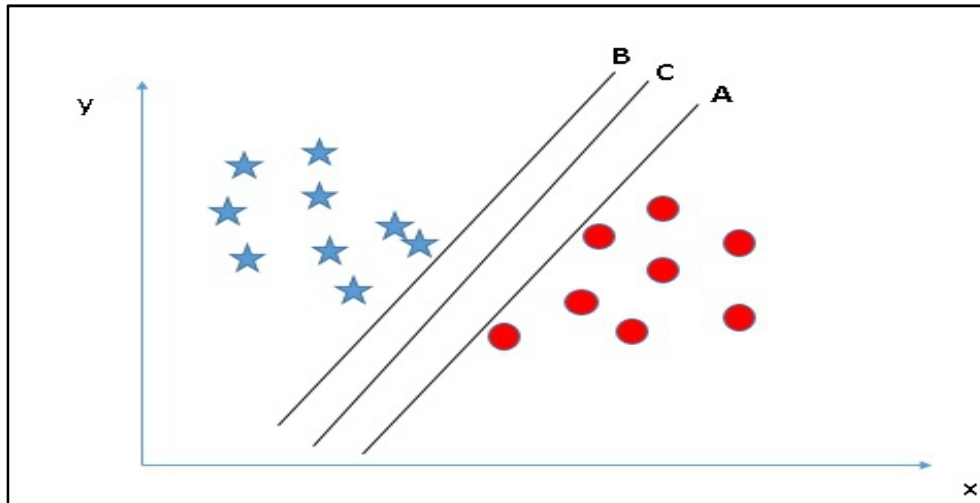[Source: https://www.analyticsvidhya.com/blog/2017/09/understaing-support-vector-machine-example-code/]

Figure 8.5(b): Hyperplane

[Source: https://www.analyticsvidhya.com/blog/2017/09/understaing-support-vector-machine-example-code/]

Hyper-plane "C" performed an excellent job in this Fig 8.5(b)

Choosing the appropriate hyper-plane will be helped by maximizing the distances between the closest data point (of any class) and the hyper-plane. We refer to this gap as a margin.
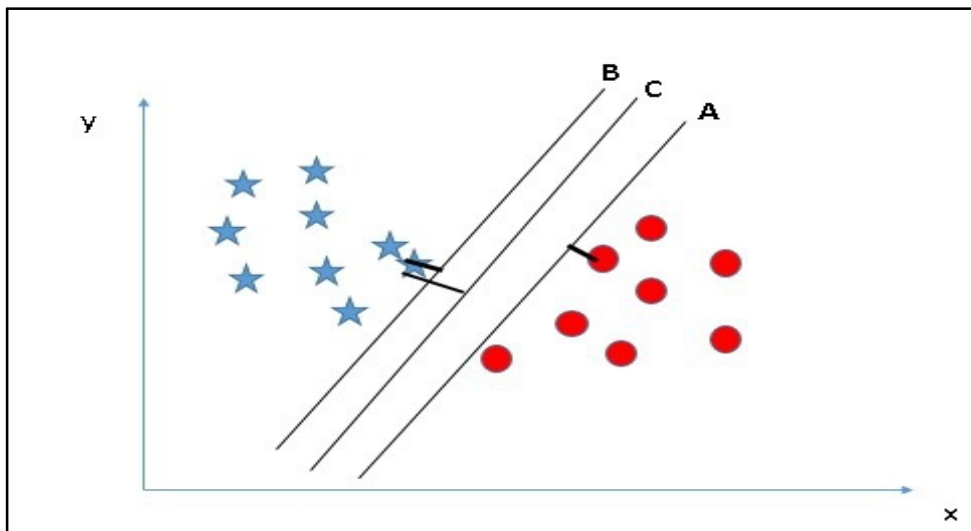


Figure 8.5(c): Hyperplane

[Source: https://www.analyticsvidhya.com/blog/2017/09/understaing-support-vector-machine-example-code/]

152

As you can see in Fig 8.5(c), hyper-plane C has a larger margin than both A and B. Thus, hyperplane C is the correct hyperplane. The data points nearest to the hyperplane are known as support vectors as shown in Fig 8.6 (a) and Fig 8.6 (b) . These points are important because they define the hyperplane's orientation and position. The hyperplane's location may shift if a support vector is deleted.
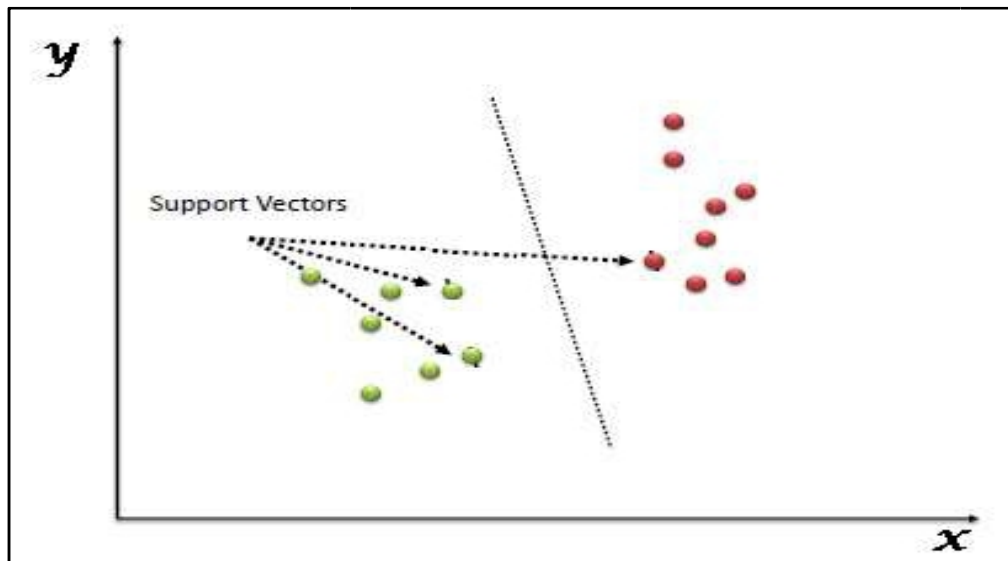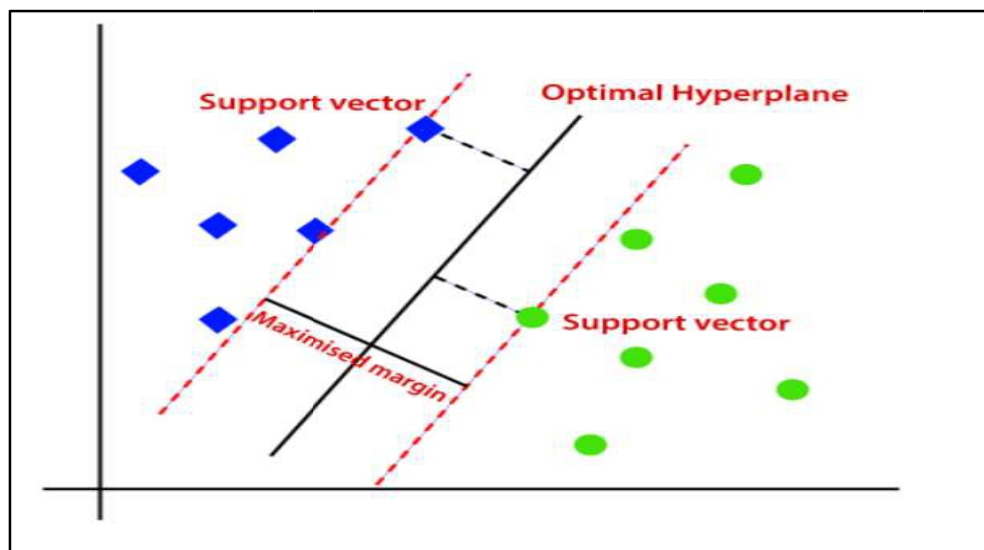


Figure 8.6(a) : Support Vectors

[Sourcs: https://www.analyticsvidhya.com/blog/2017/09/understaing-support-vector-machine-example-code/]



Figure 8.6 (b) : Support Vectors

[Source: https://www.javatpoint.com/machine-learning-support-vector-machine-algorithm]

### 8.7.1 Types of Support Vector Machines:

There are two types of SVMs as described here.

*Linear SVM*

When the data can be divided into classes using a straight line in two dimensions or a flat plane in three dimensions, then this type of SVM is known as the linear support vector machine (SVM). The optimal hyperplane for classifying the data is identified using the SVM algorithm.

*Non-Linear SVM*

When the data cannot be separated linearly, then non-linear SVM is used. In these situations, SVM uses kernel functions to move the data into a higher-dimensional space where linear separation is feasible. The best hyperplane in this new space is then found by the algorithm. Preprocessing techniques are used to convert the training data into a higher-dimensional feature space so that it may be linearly separated. However, larger dimensional spaces can become computationally demanding and raise the risk of overfitting the data, which can lead to increased complexity. The "kernel trick" substitutes a similar kernel function for the dot product computations, thus minimizing the amount of that complexity and increasing processing efficiency.

### 8.7.2 Kernel Functions in SVM

Because SVM is capable of handling both linear and nonlinear classification tasks, the SVM method is widely used in machine learning. Kernel functions are used to transform the data into a higher-dimensional space to enable linear separation, even particularly when the data cannot be separated linearly. This usage of kernel functions is sometimes referred to as the "kernel trick." The particular use case and the properties of the data determine the kernel function to apply.

Or we can say that Low-dimensional input space can be converted into a higher-dimensional space using kernel functions. Using kernel functions, SVM can generate complicated decision boundaries. These are a few well-known kernel functions mentioned in Table 8.1. The particular use case and the properties of the data determine which kernel function is to apply.

**Table 8.1:** Kernel of SVM

| Name | Formula | Description |
|------|---------|-------------|
| **Linear** | $K(x_i, x_j) = x_i \cdot x_j$ | Used when the data is linearly separable. |
| **Polynomial** | $K(x_i, x_j) = (x_i \cdot x_j + c)^d$ | Where c is a constant, and d is the degree of the polynomial. This kernel is useful for classifying data with polynomial relationships. |
| ***Radial Basis Function (RBF) Kernel / Gaussian Kernel*** | $K(x_i, x_j) = \exp(-\boldsymbol{\gamma} \lVert x_i - x_j \rVert^2)$ | Where $\gamma$ is a parameter that defines the influence of a single training example. This is one of the most popular kernels for non-linear data. |
| ***Sigmoid Kernel*** | $K(x_i, x_j) = \tanh(\boldsymbol{\alpha}\, x_i \cdot x_j + c)$ | Where α and c are kernel parameters. It behaves like a neural network's activation function. |

**8.7.3 Advantages of SVM**

**Effective in High Dimensions:** Because SVM aims to maximize the margin, it performs best when there are more features than samples.

**Robust to Overfitting:** Support Vector Machines (SVMs) are less likely to overfit when they concentrate on optimizing the margin, particularly in high-dimensional domains. **Versatile:** By using various kernels, SVM may be applied to both linear and non-linear classification tasks.

**8.7.4 Challenges**

**Selection of Kernel:** The parameters and kernel selection have a significant impact on SVM performance. Good performance depends on choosing the appropriate kernel and fine-tuning its settings. **Computationally Expensive:** Because SVM solves a quadratic optimization issue, it can be computationally costly, particularly for big datasets.

**Sensitive to Outliers:** Since outliers have the potential to alter the hyperplane's location and lower the margin, SVM may be inclined to them.

### 8.7.5 Applications of SVM

There are several uses for SVM, such as:

**Text and Hypertext Classification:** Arranging documents in groups according to specific criteria, such as spam.

**Image classification:** Identifying items in images or grouping pictures according to their contents.
**Bioinformatics:** Gene and protein classification, illness prognostication.
**Handwriting Recognition:** Categorizing letters or numbers written by hand.

---

*Check Your Progress 6.*

Q1. Define hyperplane.

Q2. What is support vector in SVM?

Q3. What is role of kernel in SVM?

Q4. What are the major challenges of SVM?

Q.5 Write names of major kernel functions.

---

## 8.8 Summary

This chapter covers the basic concept of Analytics, Classification, and prediction its need, and its spectrum. The chapter also focuses on learning and using essential approaches for making intelligent forecasts from historical data. It includes major approaches: classification, which categorizes data into specified groups (for example, spam detection), and prediction. A vast variety of models and techniques are used to handle both Classification and Prediction issues of predictive analytics. It is important to select the appropriate methodology based on the situation, the type of data, and the available computer resources, as each method has pros and cons of its own. While models like KNN and Bayesian classifiers are renowned for their simplicity and interpretability, techniques like decision trees, neural networks, and ensemble approaches are famous for their efficacy. While SVM is effective for high-dimensional data and works well with both linear and non-linear separations.

# 8.9 Terminal Questions

1. Explain the difference between classification and prediction in the context of predictive analytics.

2. What are the key challenges in predictive modeling, and how can they be addressed?

3. Explain the role of machine learning in advanced analytics. How does it enhance decision-making processes?

4. What is the role of the kernel function in an SVM? Differentiate different types of Kernel used in SVM and the cases where they are used.

5. Compare the performance of a k-nearest neighbors classifier and a support vector machine on the same dataset. Discuss which performed better and why.

6. Develop a predictive model to forecast the likelihood of customer churn for a telecom company. What features would you consider, and how would you validate your model?

7. Imagine you are tasked with predicting the probability of loan default for a bank. What data would you need, and which classification technique would you choose? Justify your choice.

8. Explain the concept of margin in SVM. Why is maximizing the margin important?

# UNIT-9  MINING FREQUENT PATTERNS, ASSOCIATIONS AND CORRELATIONS

**Structure**

## 9.0   Introduction

The task of uncovering patterns and associations within large datasets is at the heart of data mining and knowledge discovery. There are number of techniques and methods employed in the field of data mining. Among them the identification of frequent patterns stands out as a fundamental and powerful tool. Mining frequent patterns is a foundational technique in data mining and is crucial for uncovering hidden insights and making informed decisions across various domains. Frequent pattern mining is essential not only for understanding the underlying structure and relationships within data but also for enabling a wide range of applications, from market basket analysis to bioinformatics. These patterns appear frequently and can reveal important relationships and associations in the data. This process involves relevant algorithms and techniques to ensure efficient and accurate discovery of valuable insights.

Frequent patterns, including item-sets, sub-sequences, and substructures, are patterns that appear frequently within a dataset. The discovery of these patterns is essential for various applications. For instance, in market basket analysis, frequent item-sets reveal products that are often purchased together, providing invaluable insights for inventory management, marketing strategies, and customer relationship management. Similarly, in bioinformatics, frequent patterns can help identify gene sequences associated with specific diseases, aiding in the advancement of medical research and personalized medicine.

The importance of frequent pattern mining is underscored by its foundational role in several advanced data mining tasks. Techniques such as association rule learning, sequence mining, and web mining build upon the principles of frequent pattern mining, demonstrating its versatility and significance. This chapter will explore the theoretical concepts of frequent pattern mining, examine various algorithms designed to efficiently uncover these patterns, and discuss their applications in real-world scenarios. Understanding frequent pattern mining requires a details of several key concepts and challenges. These include the definition of support and confidence, the handling of large and sparse datasets, and the development of scalable and efficient algorithms. This unit is equipped with the knowledge and tools necessary to effectively mine frequent patterns, defining the way for deeper insights and more informed decision-making in their respective fields.

## 9.1 Objectives

Objectives of this unit are:

a) To introduce the Core Concepts of Mining frequent patterns.

b) To describe a number of different Techniques and Algorithms.

c) To describe Multilevel and Multidimensional Analysis used in mining frequent patterns

d) To enabling the discovery of valuable insights that can inform decision-making, improve processes, and enhance understanding of data in various domains.

## 9.2    Basic Concepts

Mining frequent patterns is a crucial step in data mining, particularly in association rule learning and market basket analysis. It involves discovering patterns, associations, correlations, or causal structures among sets of items in transactional databases. Here are some basic concepts and steps involved in mining frequent patterns:

### 9.2.1 Item-set

An item-set is a collection of one or more items. For example, in a supermarket, an itemset could be {bread, milk, butter}.

### 9.2.2 Support

Support is a measure of how frequently an itemset appears in the dataset. It is defined as the proportion of transactions in the database in which the itemset appears.

$$\text{Support}(X) = \frac{\text{Number of transactions containing } X}{\text{Total number of transactions}}$$

### 9.2.3 Frequent Item-set

An item-set is considered frequent if its support is greater than or equal to a user-defined minimum support threshold.

### 9.2.4 Association Rule

An association rule is an implication of the form A⇒B, where A and B are itemsets. The rule suggests that transactions containing itemset A are likely to contain itemset B.

### 9.2.5 Confidence

Confidence is a measure of the reliability of an association rule. It is defined as the proportion of transactions containing AAA that also contain BBB.

$$\text{Confidence}(A \Rightarrow B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A)}$$

### 9.2.6 Lift

Lift is a measure of how much more likely B is to occur when A has occurred compared to when A has not occurred. It is defined as:

$$\text{Lift}(A \Rightarrow B) = \frac{\text{Confidence}(A \Rightarrow B)}{\text{Support}(B)}$$

### 9.2.7 Steps in Mining Frequent Patterns

1. **Generate Candidate Itemsets**
   - Start with itemsets of length 1.

- Generate candidate itemsets of length k from frequent itemsets of length k−1.

2. **Prune Infrequent Itemsets**
   - Remove candidate itemsets that do not meet the minimum support threshold.

3. **Calculate Support**
   - Scan the database to determine the support of each candidate itemset.

4. **Determine Frequent Itemsets**
   - Identify itemsets that meet the minimum support threshold.

5. **Generate Association Rules**
   - For each frequent itemset, generate all possible association rules and calculate their confidence.
   - Retain rules that meet the minimum confidence threshold.

**Applications of Frequent Pattern Mining**

- **Market Basket Analysis:** Understanding buying behavior by finding associations between products.
- **Intrusion Detection:** Identifying patterns of activities that suggest security breaches.
- **Web Usage Mining:** Discovering user navigation patterns on websites.
- **Bioinformatics:** Finding patterns in biological data, such as DNA sequences.

Understanding and applying these concepts allows for extracting meaningful insights from large datasets, driving decisions in various domains such as retail, finance, and healthcare.

---

*Check Your Progress 1:*

Q.1 What do you mean by mining the frequent patterns?

Q.2 What are the core concepts of Mining the frequent patterns?

Q.3 Differentiate between the support and Confidence value?

Q.4 Define the basic Steps used for Mining Frequent Patterns?

Q.5 How the Frequent pattern mining helps in different sectors?

---

## 9.3 Frequent Item-Set Mining Algorithms

Frequent item-set mining algorithms are essential for discovering interesting patterns in large transactional databases. Here are some of the most well-known algorithms:

### 9.3.1 Apriori Algorithm

The Apriori algorithm is one of the earliest and most well-known algorithms for frequent itemset mining.

**Key Concepts:**

- **Join Step:** Combine frequent itemsets of size k to generate candidate itemsets of size k+1.
- **Prune Step:** Remove candidate itemsets that have infrequent subsets.

**Steps:**

1. Identify all frequent 1-itemsets by scanning the database.
2. Generate candidate 2-itemsets from frequent 1-itemsets.
3. Scan the database to find the support of each candidate 2-itemset and prune those that do not meet the minimum support threshold.
4. Repeat the process for k-itemsets until no more frequent itemsets are found.

**Pros:** Simple and easy to implement.

**Cons:** Can be inefficient due to the generation of a large number of candidate itemsets and multiple database scans.

**Example:**

Suppose we have the following transactional database with a minimum support threshold of 50% (i.e., an itemset must appear in at least 2 out of 4 transactions to be considered frequent):

| Transaction ID | Items Purchased |
|---|---|
| T1 | A, B, C |
| T2 | A, C |
| T3 | A, D |

| Transaction ID | Items Purchased |
|---|---|
| T4 | B, C |

**Steps:**

1. **Generate Candidate Itemsets (C1):**

   - Count the occurrence of each item in all transactions.

   C1 = {A: 3, B: 2, C: 3, D: 1}

2. **Generate Frequent Itemsets (L1):**

   - Keep only those items that meet the minimum support threshold (50% or 2 transactions).

   L1 = {A: 3, B: 2, C: 3}

3. **Generate Candidate Itemsets (C2):**

   - Generate all possible pairs of items from L1 and count their occurrences.

   C2 = {AB: 1, AC: 2, BC: 2}

4. **Generate Frequent Itemsets (L2):**

   - Again, retain only those item pairs that meet the minimum support threshold.

   L2 = {AC: 2, BC: 2}

5. **Generate Candidate Itemsets (C3):**

   - Generate candidate itemsets of size 3 from L2. Here, there's only one candidate: ABC. However, its support is 1, which is below the threshold, so we stop here.

6. **Generate Association Rules:**

   - From the frequent itemsets, generate rules like:
     - "If A, then C" (Support = 2, Confidence = 2/3 = 66.7%).

**9.3.2. FP-Growth Algorithm (Frequent Pattern Growth)**

FP-Growth is an efficient and scalable method for mining frequent itemsets without candidate generation.

**Key Concepts:**

- **FP-Tree (Frequent Pattern Tree):** A compact data structure that stores the database in a compressed form.
- **Recursive Mining:** Uses a divide-and-conquer approach to recursively find frequent itemsets.

**Steps:**

1. Scan the database to find the support of each item and sort items by descending support.
2. Construct an FP-tree by inserting transactions in the order of sorted items.
3. Mine the FP-tree recursively to extract frequent itemsets.

**Pros:** More efficient than Apriori for large datasets due to the reduced number of database scans and the compact FP-tree structure.

**Cons:** The FP-tree can become very large for dense datasets.

**Example:** Using the same transactional database:

**Steps:**

1. **Construct the FP-Tree:**
   - First, count the frequency of each item.

     {A: 3, B: 2, C: 3, D: 1}

   - Then, order the items in each transaction by descending frequency and build the FP-Tree.

     T1: A, C, B

     T2: A, C

     T3: A, D

     T4: C, B

FP-Tree construction might look like this:

(Root)

```
├─A:3

│  └─C:2

│     └─B:1

└─C:1

   └─B:1
```

2. **Mine Frequent Patterns from FP-Tree:**

- Start from the bottom (e.g., B), create conditional pattern bases, and recursively generate frequent patterns.
- From this tree, we might derive frequent patterns like (C), (A, C), (A, B, C).

## 9.3.3. **Eclat Algorithm (Equivalence Class Clustering and bottom-up Lattice Traversal)**

Eclat is a depth-first search algorithm that uses a vertical data format.

**Key Concepts:**

- **Tidsets:** Each item is associated with a list of transaction IDs (tids) in which it appears.

- **Intersection:** Frequent itemsets are found by intersecting tidsets.

**Steps:**

1. Transform the database into a vertical data format where each item is associated with a tidset.

2. Use depth-first search to explore itemsets by intersecting tidsets to count support.

3. Prune itemsets that do not meet the minimum support threshold.

**Pros:** Efficient for datasets with many transactions but fewer items.

**Cons:** Can be memory-intensive for large tidsets.

**Example:** Same transactional database:

**Steps:**

1. **Convert to Vertical Format:**

   A: {T1, T2, T3}

   B: {T1, T4}

   C: {T1, T2, T4}

   D: {T3}

2. **Intersect TID Lists to Find Frequent Itemsets:**
   - For example, to find frequent itemsets of size 2:

   A ∩ C = {T1, T2} (support = 2)

   B ∩ C = {T1, T4} (support = 2)

   - Both (A, C) and (B, C) are frequent itemsets with support ≥ 2.

9.3.4. **Direct Hashing and Pruning (DHP)**

DHP enhances the efficiency of the Apriori algorithm by reducing the number of candidate itemsets.

**Key Concepts:**

- **Hash Table:** Uses a hash table to filter out candidate itemsets that are unlikely to be frequent.

**Steps:**

1. Generate candidate 2-itemsets and store them in a hash table.
2. Use the hash table to prune infrequent itemsets early.
3. Continue generating and pruning candidate itemsets using the hash table.

**Pros:** Reduces the number of candidate itemsets.

**Cons:** Still requires multiple database scans.

**Example:** Consider the same database:

**Steps:**

1. **Generate Candidate 2-itemsets:**
   - As in Apriori, generate candidate pairs from L1.

     C2 = {AB, AC, BC}

2. **Hashing for Pruning:**
   - Use a hash function to map each candidate pair to a bucket in a hash table, and count the number of pairs in each bucket.
   - Suppose AC and BC hash to the same bucket and their combined count exceeds the minimum threshold, they are retained; otherwise, prune.

3. **Generate Frequent Itemsets:**
   - Only keep the candidate itemsets that survive the hashing and meet the minimum support.

9.3.5. **Multi-Level Mining**

Multi-level mining algorithms consider hierarchical structures in data, allowing the discovery of frequent itemsets at different levels of abstraction.

**Key Concepts:**

- **Concept Hierarchies:** Uses hierarchies (e.g., item-category-subcategory) to mine patterns at multiple levels.

**Steps:**

1. Identify frequent itemsets at a higher level of abstraction.
2. Drill down to find frequent itemsets at lower levels of abstraction.
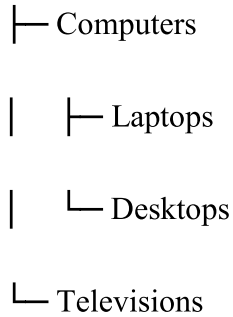3. Repeat the process for all levels of the hierarchy.

**Pros:** Provides more comprehensive insights by considering different levels of detail.

**Cons:** Complexity increases with the number of hierarchy levels.

**Example:**

Consider a retail database with a hierarchy like:

```
Electronics

├── Computers

│   ├── Laptops

│   └── Desktops

└── Televisions
```

**Steps:**

1. **Mine at Higher Level:**
   - First, mine frequent patterns at the "Electronics" level (e.g., {Electronics}).
2. **Drill Down to Lower Levels:**
   - If "Electronics" is frequent, mine at the next level (e.g., {Computers}) and continue to more specific levels (e.g., {Laptops}).
3. **Adjust Support Thresholds:**
   - Support thresholds may be different at different levels; higher-level patterns usually require higher support.

### 9.3.6. Rapid Association Rule Mining (RARM)

RARM is designed to improve the efficiency of frequent itemset mining by reducing the number of database scans and candidate generations.

**Key Concepts:**

- **Single Database Scan:** Attempts to minimize the number of database scans required.

**Steps:**

1. Perform a single scan of the database to gather initial statistics.
2. Use these statistics to generate and prune candidate itemsets efficiently.

**Pros:** Reduces computational overhead.

**Cons:** May require more sophisticated implementation techniques.

**Example:**

Let's assume we're working with a large dataset.

**Steps:**

1. **Apply Optimizations:**
   - Implement parallel processing to split the dataset and mine patterns concurrently.
   - Use data structures that minimize memory use and speed up support counting (e.g., bitmaps).

2. **Generate Frequent Itemsets and Rules Quickly:**
   - The algorithm would quickly generate frequent itemsets and then derive association rules, focusing on reducing computational overhead.

3. **Analyze Results:**
   - RARM might provide results faster than traditional methods, particularly useful in real-time applications or with very large datasets.

---

*Check Your Progress 2:*

Q.1 Explain the steps of Apriori Algorithm.

Q.2 What do you mean by FP-Growth Algorithm (Frequent Pattern Growth)?

Q.3 Define Eclat Algorithm (Equivalence Class Clustering and bottom-up Lattice Traversal).

Q.4 Discuss Direct Hashing and Pruning (DHP).

Q.5 What are steps of Multi-Level Mining?

Q.6 Define Rapid Association Rule Mining (RARM).

---

## 9.4  Mining Various Kinds Of Association

"Mining Various Kinds of Association" refers to the process of discovering relationships or patterns among a set of items, variables, or objects in large datasets. Association mining is a key concept in data mining and knowledge discovery, particularly in areas like market basket analysis, recommendation systems, and customer behaviour analysis. The term "various kinds of association" emphasizes that there are different types of associations or relationships that can be identified depending on the context and the nature of the data.

## 9.4.1 Association Rules:

- Association rules are implications of the form X⇒Y, where X and Y are disjoint itemsets (sets of items). The rule suggests that if a certain itemset X occurs, then the itemset Y is likely to occur as well.
- **Support**: The support of an itemset is the proportion of transactions in the dataset that contain the itemset. It measures how frequently the items appear together.
- **Confidence**: The confidence of an association rule X⇒Y is the probability that transaction containing X also contains Y. It measures the strength of the implication.
- **Lift**: The lift of a rule is the ratio of the observed support to that expected if X and Y were independent. Lift greater than 1 indicates a positive association, less than 1 indicates a negative association.

## 9.4.2  Types of Associations:

- **Positive and Negative Associations**: Positive associations occur when the presence of one itemset increases the likelihood of another itemset occurring. Negative associations, on the other hand, occur when the presence of one itemset decreases the likelihood of another.
- **Symmetric and Asymmetric Associations**: Symmetric associations imply that the relationship between two itemsets is bidirectional (i.e., X⇒Y and Y⇒X ). Asymmetric associations mean that the relationship is unidirectional.
- **High-Order Associations**: These involve relationships between more than two itemsets (e.g., X1, X2⇒Y ).

## 9.4.3  Mining Process:

- **Frequent Itemset Mining**: The first step in association rule mining is to identify all itemsets that occur frequently in the dataset, known as frequent itemsets. Algorithms like Apriori, Eclat, and FP-Growth are commonly used for this purpose.
- **Generating Association Rules**: Once the frequent itemsets are identified, association rules are generated based on user-defined thresholds for support and confidence.
- **Filtering and Validation**: The generated rules are filtered to remove redundant or insignificant rules, often using measures like lift, conviction, or chi-square test.

### 9.4.4 Extensions and Variants:

- **Temporal Associations**: These capture relationships that involve time. For example, a temporal association might reveal that certain products are bought together more frequently during a specific season.
- **Sequential Patterns**: These involve discovering patterns where the order of items matters, such as identifying the sequence in which products are purchased.
- **Multilevel Associations**: These consider hierarchical relationships among items. For instance, rules might be mined at different levels of abstraction, such as finding associations at the category level and then at the individual product level.
- **Multidimensional Associations**: These take into account additional dimensions like location, time, or customer demographics, allowing for more complex and informative rules.

---

*Check Your Progress 3:*

> **Q.1 What do you mean by association rule?**
> **Q.2 Define different types of associations.**
> **Q.3 Discuss the mining process of association rule.**

---

# 9.5 Rules – Multilevel And Multidimensional

**Mining frequent patterns, associations, and correlations** is a fundamental task in data mining, aimed at discovering interesting relationships between data items in large datasets. This process is often conducted through the mining of **association rules**, where patterns or correlations between data items are identified. Within this domain, two advanced concepts—**multilevel** and **multidimensional** mining—are particularly important for uncovering more nuanced and context-rich patterns.

9.5.1 Multilevel Mining

**Multilevel mining** involves discovering patterns and associations at different levels of abstraction or granularity within a hierarchical structure. This hierarchy could be based on a taxonomy or ontology, where data items are categorized into levels (e.g., categories, subcategories, items).

*Key Concepts:*

- **Hierarchy of Items**: Items can be organized in a hierarchy, such as:
  - Level 1: "Beverages"
  - Level 2: "Soft Drinks," "Juices"
  - Level 3: "Coke," "Pepsi" under "Soft Drinks"; "Orange Juice," "Apple Juice" under "Juices"

  Mining at different levels allows discovering associations that are either more general (high-level) or more specific (low-level).

- **Support Count Adjustment**: When moving between levels in a hierarchy, the support count (the frequency of an itemset) needs to be adjusted. Higher-level itemsets usually have higher support because they encompass more specific items. For example, "Soft Drinks" might have higher support than "Coke" alone.

- **Top-Down Approach**: Multilevel mining often follows a top-down approach, where patterns are first identified at a higher level and then progressively refined to lower levels. This can help in efficiently focusing on the most promising regions of the data hierarchy.

*Example:*

Consider a retail dataset where items are organized in a hierarchy:

- **High-Level Rule**: "Customers who buy beverages also tend to buy snacks."
- **Low-Level Rule**: "Customers who buy soft drinks also tend to buy chips."

Here, the high-level rule provides a broader insight, while the low-level rule gives a more specific association. Mining at multiple levels allows businesses to make decisions that are both general (for overall strategy) and specific (for targeted promotions).

*Applications:*

- **Retail and Market Basket Analysis**: Understanding purchasing patterns at various product levels (e.g., product categories, brands).
- **Customer Segmentation**: Analyzing customer behaviour at different levels of detail (e.g., demographic categories, purchasing habits).

- **Healthcare**: Discovering associations between symptoms, diagnoses, and treatments at different levels (e.g., disease categories, specific conditions).

9.5.2 Multidimensional Mining

**Multidimensional mining** extends the traditional association mining by considering multiple dimensions of the data simultaneously. This approach is useful when data items are associated with multiple attributes (or dimensions), such as time, location, customer demographics, and product categories.

*Key Concepts:*

- **Dimensions and Attributes**: Each dimension represents an aspect of the data, such as:
    - **Time**: Day, month, season
    - **Location**: Country, city, store
    - **Product**: Category, brand, item
    - **Customer**: Age, gender, income level
- **Mining Multidimensional Rules**: Instead of focusing on a single dimension (e.g., products), multidimensional mining discovers patterns across multiple dimensions. This could involve mining rules like "Customers of age 20-30 in New York tend to buy sports shoes during the holiday season."
- **Data Cubes**: Multidimensional data is often represented in the form of data cubes, where each cell represents an aggregation of data (e.g., total sales) for a particular combination of dimension values (e.g., sales of sports shoes in New York in December).
- **Constraint-Based Mining**: Users can apply constraints to focus on specific dimensions or combinations of dimensions that are of interest, making the mining process more efficient and targeted.

*Example:*

Consider a sales dataset with dimensions such as time (month), location (store), product (category), and customer (age group):

- **Multidimensional Rule**: "In December, young adults (age 20-30) in urban areas tend to buy electronics, particularly smartphones."

This rule incorporates multiple dimensions—time (December), location (urban areas), customer demographics (age 20-30), and product type (electronics, smartphones).

*Applications:*

- **Targeted Marketing**: Identifying customer segments that are likely to respond to specific promotions based on multiple factors.
- **Supply Chain Management**: Understanding how different factors like time, location, and product type affect demand.
- **Healthcare**: Analyzing patient data across multiple dimensions, such as age, gender, medical history, and geographic location, to identify patterns in disease prevalence or treatment outcomes.

---

*Check Your Progress 4:*

Q.1 Differentiate between Multi-Level association and multidimensional association.

Q.2 Discuss the application of Multi-Level association.

Q.3 Discuss the application of Multidimensional association.

---

# 9.6. Association Rule Mining Vs Correlation Analysis

**Association Rule Mining** and **Correlation Analysis** are both techniques used in data mining and statistics to uncover relationships between variables or items in large datasets. While they share similarities, they are distinct in their goals, methodologies, and applications.

**Purpose**: Correlation analysis aims to measure the strength and direction of the linear relationship between two numerical variables. Unlike association rule mining, which deals with finding patterns in categorical data, correlation analysis typically deals with continuous data.

### 9.6.1 Comparison: Association Rule Mining vs. Correlation Analysis

| Feature | Association Rule Mining | Correlation Analysis |
|---|---|---|
| Data Type | Categorical (e.g., transactions, discrete items) | Continuous (e.g., numerical data, scores) |
| Objective | Discover patterns and rules (e.g., item co- | Measure strength and direction of linear |

| Feature | Association Rule Mining | Correlation Analysis |
|---|---|---|
| | occurrences) | relationships |
| Output | Association rules (e.g., X⇒Y) | Correlation coefficient (e.g., Pearson r) |
| Measures | Support, Confidence, Lift | Pearson r, Spearman ρ, Kendall τ |
| Interpretation | Provides actionable rules, e.g., items frequently bought together | Provides a numerical value representing the relationship |
| Applicability | Market basket analysis, recommendation systems, fraud detection | Economics, healthcare, social sciences, any field involving continuous data |
| Dependency | Assumes no linearity; looks for associations based on frequency and co-occurrence | Assumes linearity (in the case of Pearson) |
| Scalability | Can handle large transactional datasets with many items | Typically applied to smaller datasets with fewer variables |

## 9.6.2 Key Differences

1. **Nature of Relationships**:
   - **Association Rule Mining**: Focuses on finding if one item or set of items is associated with another. It is not concerned with the directionality or linearity but rather with the occurrence patterns in the data.
   - **Correlation Analysis**: Specifically measures how two variables move together in a linear fashion, indicating the direction (positive or negative) and strength of the relationship.
2. **Data Handling**:
   - **Association Rule Mining**: Works well with categorical data, particularly in large datasets like market baskets. It can handle multiple items in a transaction and finds associations across these items.
   - **Correlation Analysis**: Primarily deals with continuous, numerical data and is more about understanding the relationship between two specific variables rather than finding patterns across many variables.

3. **Interpretation of Results**:
   - **Association Rule Mining**: Results are actionable rules that can be directly applied in decision-making processes, like promoting items that are frequently bought together.
   - **Correlation Analysis**: Results are more about understanding the relationship between variables and are used in modeling, prediction, and hypothesis testing.

### 9.6.3 When to Use Which?

- **Use Association Rule Mining** when you need to discover patterns or associations in categorical data, particularly in transactional or categorical datasets. It's useful in scenarios like market basket analysis, where you want to know which items are often purchased together.
- **Use Correlation Analysis** when you need to understand the linear relationship between two continuous variables. This is applicable in scenarios where you are interested in how one variable might predict or explain the changes in another, such as in regression analysis or when testing hypotheses about relationships between variables.

---

*Check Your Progress 5:*

**Q.1 Differentiate between Association Rule Mining vs. Correlation Analysis.**

**Q.2 Discuss the application of Multi-Level association.**

**Q.3 Discuss the application of Multidimensional association.**

---

# 9.7 Summary

These algorithms provide various approaches to mining frequent itemsets, each with its own strengths and weaknesses. The choice of algorithm depends on the specific characteristics of the dataset and the requirements of the analysis task. **Apriori Algorithm** is a classic algorithm for mining frequent itemsets and generating association rules, based on iterative candidate generation and support counting. **FP-Growth Algorithm** is an efficient alternative to Apriori that uses a compact FP-Tree structure to mine frequent patterns without candidate generation. **Eclat Algorithm** is avertical data format-based algorithm that efficiently finds frequent itemsets by intersecting transaction ID lists. **Direct Hashing and Pruning (DHP)**is an optimization technique that uses hash tables to prune candidate itemsets early in the mining process. **Multi-Level Mining** is a technique for discovering frequent patterns at multiple levels of

abstraction in hierarchical data. **Rapid Association Rule Mining (RARM)** is a method focused on speeding up the process of mining association rules, suitable for large-scale datasets.

**Multilevel and multidimensional mining** of frequent patterns, associations, and correlations provides a richer understanding of the data by considering different levels of abstraction and multiple attributes simultaneously. These techniques are powerful tools in various fields, including retail, marketing, healthcare, and beyond, enabling more informed decision-making through the discovery of complex, context-aware patterns. While both association rule mining and correlation analysis aim to uncover relationships in data, they differ significantly in their methods, types of data they analyze, and the kinds of relationships they reveal. Understanding these differences is crucial for applying the correct technique to the appropriate problem, leading to more effective data analysis and decision-making.

# 9.8 Terminal Questions

Q1. What is an example of an association rule in mining?

Q2. If {1, 2, 3}, {1, 2, 4}, {1, 3, 4}, {1, 2, 5}, and {3, 4, 5} are ALL the large 3-itemsets, list all of the large 2-itemsets out of the original data set - {1, 2, 3, 4, 5}.

Q3. Consider the Data set D. Given the minimum support 2, apply Apriori algorithm on this dataset.

| Transaction ID | Items |
|---|---|
| 100 | A,C,D |
| 200 | B,C,E |
| 300 | A,B,C,E |
| 400 | B,E |

Q4. Write and explain the algorithm for mining frequent item sets candidate generation. Give relevant example.

Q.5 Discuss the approaches for mining multi-level association rules from the transactional databases. Give relevant example.

**Master of Computer Application**

# MCA-E6N
# Data Mining

उ० प्र० राजर्षि टण्डन
मुक्त विश्वविद्यालय, प्रयागराज

# Block
# 4

## Advanced Data Mining Techniques

## Course Design Committee

| | |
|---|---|
| **Prof. Ashutosh Gupta** | **Chairman** |
| Director (In-charge) | |
| School of Computer & Information Science, UPRTOU Allahabad | |
| **Prof. Suneeta Agarwal** | **Member** |
| Dept. of Computer Science & Engineering | |
| Motilal Nehru National Institute of Technology Allahabad | |
| **Dr. Upendra Nath Tripathi** | **Member** |
| Associate Professor | |
| DeenDayalUpadhyay Gorakhpur University, Gorakhpur | |
| **Dr. Ashish Khare** | **Member** |
| Associate Professor | |
| Dept. of Computer Science, University of Allahabad, Prayagraj | |
| **Ms. Marisha** | **Member** |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |
| **Mr. Manoj Kumar Balwant** | **Member** |
| Assistant Professor (Computer Science) | |
| School of Science, UPRTOU Allahabad | |

## Course Preparation Committee

| | |
|---|---|
| **Dr. Tulika Narang** | **(Block 1, 2 & 4) Author** |
| Assistant Professor, Computer science | |
| United University, Rawatpur, Prayagraj | |
| **Dr. Krishan Kumar** | **(Block 3-Unit 7) Author** |
| Assistant Professor, | |
| Department of Computer Science, Faculty of Technology | |
| Gurukula Kangri Vishwavidyalaya, Haridwar (UK) | |
| **Dr. Pooja Yadav** | **(Block 3-Unit 8) Author** |
| Assistant, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Iram Naim** | **(Block 3-Unit 9) Author** |
| Assistant Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Dr. Brajesh Kumar** | **Editor** |
| Associate Professor, Departement of CS & IT | |
| Mahatma Jyotiba Phule Rohilkhand University- Bareilly | |
| **Mr. Manoj Kumar Balwant** | **Coordinator** |
| Assistant Professor (Computer Science), | |
| School of Sciences, UPRTOU, Prayagraj | |

# Block 4  Introduction

This block deeply explores web mining, a crucial step in extracting valuable knowledge from the huge data available on WWW. Web Mining is one of the least explored areas in data mining; it involves applying data mining techniques to discover patterns, trends, and natural laws from different sources, i.e., the World Wide Web. In tenth unit, we categorize the technologies into Web Content Mining, Web Structure Mining (link mining), and Web Usage Mining. This unit also compares web mining to traditional data mining and discusses the importance and use of each type. It also discusses one of the most significant concepts called Page Rank, an algorithm Larry Page and Sergey Brin created that helps determine where websites should be ranked in search engine results.

Unit eleven shifts focus to Text Mining (text data mining or text analytics). This process requires us to leverage natural language processing (NLP) techniques, machine learning, and statistics to extract the information needed from unstructured text data. The most important goal is to take plain text as input and get structured data that can then be analysed to detect patterns and trends. It also involves introducing Episode Rule Discovery and contains units on Hierarchical Categorization.

The final twelfth unit is on Spatial Mining, i.e. the process of extracting knowledge and patterns from spatial data. This unit introduces fundamental methods and algorithms of spatial mining, which focus on analysing datasets containing geographical (ex, maps and GPS coordinates) information. This process allows the spatial mining of traditional data mining techniques to incorporate and fully understand data within its geographical context.

Together, these units provide a comprehensive overview of advanced data mining techniques applicable to web, text, and spatial data, equipping students with the knowledge to leverage these methodologies in real-world scenarios.

# UNIT 10: Web Mining

**Structure**

## 10.0 Introduction

The unit focuses on the Web Mining as an essential process of finding relevant and useful knowledge from World Wide Web. Web Web mining is the process of extracting useful information from the vast amount of data available on the web. It involves applying data mining techniques to web data to discover patterns, trends, and relationships that can provide valuable insights. Web mining can be broadly divided into three main categories: The different types of web mining are also discussed in this unit. Web mining is categorised as Web Content Mining, Web structure mining and Web usage mining. Web Mining and Data mining are compared in this unit. The significance and various applications of web content mining, web structure mining and web usage mining are also discussed in this unit. Page Rank algorithm is explained in this unit with an example. PageRank is an algorithm used by search engines, most notably Google, to rank web pages in their search engine results. It was developed by Larry Page and Sergey Brin at Stanford University in the late 1990s as part of their research on a new kind of search engine.

## 10.1 Objectives

After the end of this unit, you should be able to:

- Explain Web Mining

- Explain different types of Web Mining

- Understand Web Content Mining

- Understand Web structure mining

- Understand Web usage mining

## 10.2 Web Mining

Web mining is the process of extracting useful information or knowledge from web data. It involves collecting and analyzing data from websites, web pages, web documents, and other web resources to discover patterns, trends, and insights. Web mining can be classified into three main types:
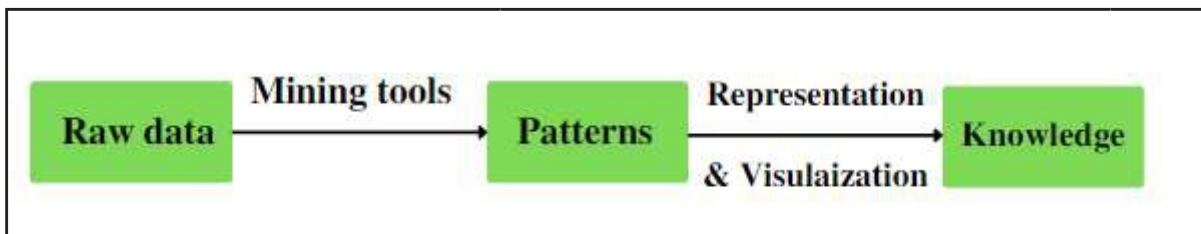


Figure1: Converting Data into Knowledge (Source: https://www.geeksforgeeks.org/web-mining/)

1. Web Content Mining: This type of web mining focuses on extracting information from the actual content of web pages. It involves techniques such as information extraction, text mining, natural language processing, and sentiment analysis to analyze and understand the textual data present on web pages.

2. Web Structure Mining: Web structure mining aims to discover the structure and organization of web pages and websites. It involves analyzing the hyperlinks between web pages, identifying patterns in linkages, and determining the relationships between different websites. This type of mining helps in tasks such as web page ranking, link analysis, and identifying communities or clusters of related websites.

3. Web Usage Mining: Web usage mining deals with analyzing user interactions and behaviour on websites. It involves tracking and analyzing data related to user visits, clicks, navigation patterns, and other actions performed by users while browsing the web. This information can be used to understand user preferences, improve website usability, personalize recommendations, and support various marketing and business decisions.
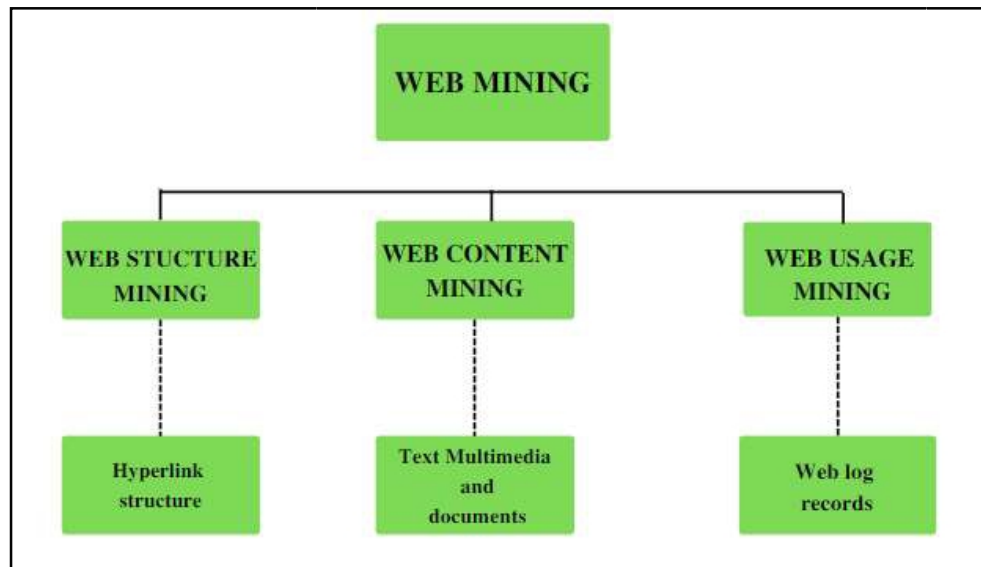
184

Figure 2: Web Mining Classification (Source: https://www.geeksforgeeks.org/web-mining/)

Web mining and data mining are related fields that involve extracting insights and knowledge from data. While they share some similarities, there are also distinct differences between the two. Data mining is a broader term that encompasses the process of discovering patterns, relationships, and insights from structured or unstructured data. It involves analyzing large datasets to uncover hidden patterns, extract meaningful information, and make predictions or decisions based on the data. Data mining techniques can be applied to various domains, including finance, healthcare, marketing, and more. On the other hand, web mining specifically focuses on extracting knowledge from web-related data, including web pages, websites, web server logs, user behaviour, and other web resources. Web mining techniques are designed to handle the unique characteristics of the web, such as the unstructured nature of web content, the presence of hyperlinks, and the dynamics of web data. Web mining can be seen as a subfield or application area of data mining that deals with data specific to the web. It combines elements of data mining, information retrieval, and machine learning to analyze web data and derive insights. Web mining techniques include web content mining, web structure mining, and web usage mining, as mentioned earlier.

While data mining can be applied to various types of data sources, including databases, text documents, and sensor data, web mining focuses specifically on data from the web. Web mining techniques often employ specialized algorithms and approaches to handle the unique characteristics and challenges posed by web data, such as the high dimensionality, noise, and the dynamic nature of the web. Data mining is a broader field that encompasses the analysis of various types of data, while web mining is a specialized area within data mining that focuses on extracting knowledge from web-related data. Web

mining techniques are tailored to handle the specific challenges and characteristics of web data. Web mining techniques often involve data pre-processing, data cleaning, feature extraction, and applying machine learning or statistical algorithms to extract meaningful insights from the collected web data. The insights gained from web mining can be applied in various domains, such as e-commerce, marketing, information retrieval, customer behavior analysis, and personalized services. Web mining faces several challenges and issues that researchers and practitioners need to address. Some of the key issues in web mining:

1. Data Quality: Web data is often noisy, unstructured, and inconsistent. Web mining algorithms need to deal with issues such as missing data, incomplete or inaccurate information, and inconsistencies in data representation. Data cleaning and pre-processing techniques are crucial to ensure the quality and reliability of mined information.

2. Scalability: The web is vast and constantly growing, making scalability a significant challenge in web mining. Processing and analyzing large volumes of web data in a timely manner require efficient algorithms and distributed computing techniques. Handling big data in web mining is crucial to extract meaningful insights without being overwhelmed by the sheer volume of information.

3. Privacy and Ethical Concerns: Web mining involves collecting and analyzing user data, which raises privacy concerns. The extraction of personal information without consent or the use of data for unintended purposes can infringe on privacy rights. It is important to ensure that web mining practices comply with privacy regulations, obtain user consent when necessary, and handle sensitive information appropriately.

4. Dynamic Nature of the Web: The web is highly dynamic, with frequent updates to web pages, new content being added, and links changing. Web mining techniques need to adapt to the evolving web environment and handle the challenges posed by changing web structures and content. Continuous monitoring and updating of web data are necessary to maintain the relevance and accuracy of mined information.

5. Semantic Gap: The web contains a vast amount of unstructured data, making it challenging to extract semantic meaning from web content. Interpreting the meaning and context of web pages accurately is a complex task. Natural language processing and machine learning techniques are employed to bridge the semantic gap and improve the understanding of web content.

6. Information Overload: The abundance of web data can lead to information overload, where users are overwhelmed by the sheer amount of available information. Web mining techniques need to focus on filtering and extracting relevant and valuable information to meet users' specific needs.

Personalization and recommendation systems can help mitigate information overload by providing tailored content and recommendations.

7. Intellectual Property Rights: Web mining involves accessing and analyzing web content, which can raise concerns regarding copyright and intellectual property rights. Web mining practitioners need to be aware of and respect the legal restrictions and intellectual property regulations associated with web data usage.

---

*Check your progress 1*

Q1. What is Web Mining?

Q2. How is data mining different from Web Mining?

Q3. Discuss the categorisation of Web Mining.

Q4. Discuss essential challenges in Web Mining.

---

## 10.3 Web Content Mining

Web content mining, also known as text mining or web text mining, is a type of web mining that focuses on extracting meaningful information and knowledge from the textual content of web pages. It involves applying various natural language processing (NLP) and text mining techniques to analyze and understand the information present in web documents. Here are some common techniques and tasks involved in web content mining:

1. Information Extraction: Information extraction techniques are used to identify and extract structured information from unstructured web content. This includes extracting entities (such as names, locations, organizations), relationships between entities, and other relevant information. Techniques like named entity recognition, entity linking, and relationship extraction are commonly used in information extraction.

2. Text Classification: Text classification is the process of categorizing web documents into predefined categories or classes based on their content. This can be achieved through supervised learning algorithms such as decision trees, support vector machines (SVM), or deep learning approaches like convolutional neural networks (CNN) or recurrent neural networks (RNN). Text classification can be used for tasks like sentiment analysis, topic classification, or content filtering.

3. Sentiment Analysis: Sentiment analysis, also known as opinion mining, aims to determine the sentiment or subjective information expressed in web content. It involves classifying the polarity of

text as positive, negative, or neutral. Sentiment analysis techniques can be used to analyze customer reviews, social media posts, or any other textual data where sentiment is expressed.

4. Text Summarization: Text summarization techniques are employed to generate concise summaries of lengthy web documents. This can be done through extractive techniques, which select and aggregate the most important sentences or phrases from the original text, or abstractive techniques, which generate new summaries by understanding the meaning and context of the content.

5. Topic Modeling: Topic modeling techniques, such as Latent Dirichlet Allocation (LDA), are used to automatically discover and extract latent topics or themes present in a collection of web documents. Topic modeling helps in organizing and categorizing web content based on underlying themes, enabling better content management and understanding of the main topics discussed.

6. Named Entity Recognition (NER): Named Entity Recognition is a technique used to identify and classify named entities, such as persons, organizations, locations, dates, or other predefined entities, in web text. NER can be used for information extraction, entity linking, or enhancing search and recommendation systems by identifying relevant entities.

Web content mining techniques are often combined with other web mining approaches, such as web structure mining and web usage mining, to provide a comprehensive understanding of web data. By extracting valuable information from web content, organizations can improve search engines, personalize user experiences, extract insights for market research, and enable various other applications that rely on understanding textual information present on the web.

# 10.4 Web Structure Mining

Web structure mining is a branch of web mining that focuses on analyzing the hyperlink structure of the web, the relationships between web pages, and the overall organization of websites. It involves extracting meaningful patterns, insights, and knowledge from the link-based structure of the web. Some essential techniques and tasks involved in web structure mining:

1. Link Analysis: Link analysis techniques examine the relationships and patterns formed by hyperlinks between web pages. This includes analyzing incoming links (backlinks) and outgoing links from a particular page, calculating link-based metrics such as PageRank, and determining the importance or authority of web pages based on their link structure.

2. PageRank Algorithm: The PageRank algorithm, developed by Google's co-founders, Larry Page and Sergey Brin, is a well-known algorithm used in web structure mining. It assigns a numerical value (PageRank score) to each web page, indicating its importance or relevance based on the number and

quality of incoming links. PageRank plays a crucial role in search engine ranking and determining the visibility of web pages in search results.

Link-Based Clustering: Link-based clustering is a method used in data analysis and machine learning to group items (such as data points or nodes in a network) based on their connections or relationships. Unlike traditional clustering techniques that rely on attributes or features of the data points themselves, link-based clustering focuses on the links (or edges) between data points. Link-based clustering techniques group related web pages based on their link patterns and similarities. Clustering algorithms, such as k-means clustering or hierarchical clustering, can be applied to identify communities or clusters of web pages that are topically or thematically similar.

3. Community Detection: Community detection techniques aim to identify groups or communities of web pages that have strong interconnectivity within themselves but are relatively less connected with other communities. Community detection helps in understanding the structural organization of the web and can be applied in tasks such as identifying related websites or analyzing social networks.

4. Web Graph Analysis: Web graph analysis involves studying the structure and properties of the web as a graph. Graph theory concepts and algorithms are employed to analyze properties such as connectivity, degree distribution, centrality measures, and graph clustering coefficients. Web graph analysis provides insights into the global structure and characteristics of the web.

5. Link Prediction: Link prediction techniques predict the likelihood of future links between web pages based on the existing link structure. These techniques can be used to discover missing or potential links, identify link spam, or suggest relevant connections in various recommendation systems. Link prediction is a task in network analysis and graph theory that involves predicting the likelihood of a future connection (or "link") between two nodes in a network based on the current structure of the network and possibly other attributes. It is widely used in various fields such as social network analysis, biology, recommender systems, and more.

Web structure mining is crucial for search engine optimization (SEO), understanding web navigation patterns, identifying authoritative sources, detecting web spam or malicious websites, and improving the overall organization and usability of websites. By analyzing the hyperlink structure, web structure mining techniques enhance our understanding of the relationships and connectivity within the web, providing valuable insights for various web-related applications.

# 10.5 Web Usage Mining

The term "web usage mining" refers to the process of analyzing and extracting insights from user behaviour and interactions on the web. It involves collecting and analyzing data related to user activities, preferences, clickstreams, and other behavioural patterns to gain a better understanding of users and their interactions with web-based systems and services. Some essential aspects and techniques associated with web user mining:

1. Web Usage Data Collection: Web user mining starts with the collection of web usage data, which can include information such as page visits, clicks, navigation paths, timestamps, session durations, and other user-related data. This data is typically collected through log files, cookies, tracking mechanisms, or analytics tools embedded in websites.

2. User Profiling: User profiling involves creating profiles of individual users based on their web usage data. This includes capturing demographic information, preferences, interests, and behavior patterns. User profiling helps in understanding user segments, personalizing content, and providing targeted recommendations.

3. Clickstream Analysis: Clickstream analysis focuses on analyzing the sequence of user clicks and navigation patterns during their web browsing sessions. It helps understand how users navigate through websites, which pages they visit most frequently, the order of their interactions, and the duration spent on each page. Clickstream analysis aids in improving website design, optimizing user experience, and identifying potential bottlenecks or usability issues.

4. Recommender Systems: Web user mining is closely related to the development of recommender systems. These systems use user behaviour data to provide personalized recommendations for products, content, or services. Collaborative filtering and content-based filtering are common techniques used in recommender systems to analyze user preferences and generate relevant recommendations.

190

5. Customer Behaviour Analysis: Web user mining enables the analysis of customer behavior, preferences, and purchasing patterns. By examining user interactions, navigation paths, and conversion rates, organizations can gain insights into customer needs, identify potential cross-selling opportunities, and make data-driven decisions to improve marketing strategies.

6. Personalization and Customization: Web user mining helps in delivering personalized experiences to users by tailoring content, recommendations, and user interfaces based on individual preferences and behaviour. Personalization techniques leverage user profiles and historical data to create customized experiences that cater to specific user needs and interests.

7. User Churn Analysis: User churn analysis focuses on understanding user attrition or churn rates in web-based services. By analyzing user behaviour patterns and identifying factors that contribute to user churn, organizations can take proactive measures to retain users, improve service quality, and enhance customer satisfaction.

Web user mining techniques require careful consideration of privacy and ethical concerns. Data collection and analysis must comply with privacy regulations and ensure user consent and anonymity when necessary. Overall, web user mining provides valuable insights into user behaviour, preferences, and interactions, enabling organizations to enhance user experiences, optimize web-based services, and make informed business decisions.

## 10.6 Web Mining Algorithms

Web mining algorithms are used to extract valuable insights and knowledge from web data. These algorithms are designed to handle the unique characteristics of web-related data, such as its unstructured nature, the presence of hyperlinks, and the dynamic nature of the web. Here are some common web mining algorithms:

1. Apriori Algorithm: The Apriori algorithm is commonly used in web mining for association rule mining. It identifies frequent itemsets in web transaction data to discover patterns of co-occurring items. This algorithm helps in understanding user behavior and preferences based on the items they interact with on web pages.

2. PageRank Algorithm: The PageRank algorithm is a well-known algorithm in web structure mining. It assigns a numerical value (PageRank score) to each web page based on the number and quality of incoming links. PageRank is a key component of search engine algorithms and helps determine the order in which search engine results are presented.

3. TF-IDF Algorithm: The Term Frequency-Inverse Document Frequency (TF-IDF) algorithm is commonly used in web content mining. It assesses the importance of words or terms within web documents. TF-IDF assigns higher weights to terms that appear frequently within a specific document but are rare in the entire collection of documents. It helps in information retrieval, text mining, and keyword extraction tasks.

4. Clustering Algorithms: Clustering algorithms, such as k-means clustering or hierarchical clustering, are used in various web mining tasks. These algorithms group similar web pages or users based on their features or behaviors. Clustering can help identify communities, detect anomalies, and provide personalized recommendations.

5. Classification Algorithms: Classification algorithms, including decision trees, support vector machines (SVM), and naive Bayes, are used in web mining for tasks such as sentiment analysis, spam detection, or content categorization. These algorithms learn from labeled data to classify web documents or user behavior into predefined categories.

6. Sequential Pattern Mining Algorithms: Sequential pattern mining algorithms, like the GSP (Generalized Sequential Pattern) algorithm, are utilized in web usage mining to discover sequential patterns in user clickstream data. They help identify the order and frequency of web page visits, allowing for the analysis of user navigation patterns.

7. Collaborative Filtering Algorithms: Collaborative filtering algorithms, such as item-based or user-based collaborative filtering, are commonly used in recommendation systems for web mining. These algorithms analyze user behaviour and preferences to make personalized recommendations, such as suggesting similar products or content based on a user's past interactions.

These are just a few examples of web mining algorithms. The choice of algorithm depends on the specific web mining task and the characteristics of the web data being analyzed. Different algorithms may be combined or adapted to address specific challenges and requirements in web mining applications. Page Rank is an algorithm used in web structure mining to determine the importance or relevance of web pages based on their link structure. Developed by Larry Page and Sergey Brin, the co-founders of Google, PageRank is a key component of Google's search engine ranking algorithm. The PageRank algorithm assigns a numerical value, known as a PageRank score, to each web page. The score is calculated based on the number and quality of incoming links to a page. A web page with more incoming links from other reputable or high-quality pages is considered more important or authoritative and will have a higher PageRank score. The calculation of PageRank involves iterative computations, where the PageRank scores of pages are updated based on the incoming links from other pages. The algorithm assumes that a page's importance is influenced not only by the quantity of links but also by the importance of the linking

pages. This approach aims to capture the idea of "voting" or "endorsement" from other pages to determine the relative significance of a page. PageRank is used by search engines to help determine the order in which search results are presented to users. Pages with higher PageRank scores are often displayed more prominently in search results. However, it is important to note that PageRank is just one of many factors that search engines consider when ranking pages, and search algorithms have evolved significantly since the original introduction of PageRank.

Since the original publication of the PageRank algorithm, variations and improvements have been developed, considering additional factors and addressing potential limitations. Nonetheless, PageRank remains a foundational algorithm in web structure mining and has greatly influenced the field of search engine optimization (SEO). A simple example to illustrate the PageRank algorithm is as follows-

Suppose we have four web pages, labeled A, B, C, and D. The following are the links between these pages:

- Page A has outgoing links to pages B and C.
- Page B has outgoing links to pages A and C.
- Page C has outgoing links to pages A and D.
- Page D has an outgoing link to page C.

To calculate the PageRank scores, we start by assuming equal initial scores for all pages. Let's say each page has an initial score of 0.25. In each iteration of the algorithm, the PageRank scores are updated based on the incoming links to each page. The updated scores are calculated using a formula that combines the current scores of linking pages and the number of outgoing links from those pages. Here's a step-by-step calculation of the PageRank scores for our example:

1. Initial scores:
   Page A: 0.25
   Page B: 0.25
   Page C: 0.25
   Page D: 0.25
2. Iteration 1:
   Page A: $(0.25 * 0.5) + (0.25 * 0.5) = 0.25$
   Page B: $(0.25 * 0.5) = 0.125$
   Page C: $(0.25 * 0.5) + (0.25 * 0.5) = 0.25$

Page D: $(0.25 * 0.5) = 0.125$

3. Iteration 2:

Page A: $(0.25 * 0.5) + (0.125 * 0.5) = 0.1875$

Page B: $(0.25 * 0.5) + (0.25 * 0.5) = 0.25$

Page C: $(0.25 * 0.5) + (0.125 * 0.5) = 0.1875$

Page D: $(0.25 * 0.5) + (0.25 * 0.5) = 0.25$

4. Iteration 3:

Page A: $(0.25 * 0.25) + (0.1875 * 0.5) = 0.171875$

Page B: $(0.25 * 0.25) + (0.25 * 0.5) = 0.1875$

Page C: $(0.25 * 0.25) + (0.1875 * 0.5) = 0.171875$

Page D: $(0.25 * 0.25) + (0.25 * 0.5) = 0.1875$

5. Iteration 4 (convergence):

Page A: 0.171875

Page B: 0.1875

Page C: 0.171875

Page D: 0.1875

After the fourth iteration, the PageRank scores have converged, indicating the final estimated importance of each page. In this example, Page B has the highest PageRank score, followed by Pages A and D, while Page C has the lowest score. This example demonstrates the basic calculation process of PageRank for a small network of web pages. In real-world scenarios, the algorithm considers a much larger number of pages and incorporates additional factors for more accurate and meaningful rankings

---

***Check your Progress 3***

Q1. What is Web Usage Mining?

Q2. Discuss some essential aspects and techniques associated with web usage mining.

Q3. Explain Page Rank algorithm.

Q4. What is the significance of Clustering and Classification in Web Mining?

---

## 10.7 Summary

The unit discussed Web Mining as an essential process of finding significant and valuable knowledge from huge repository of World Wide Web. The process of Web mining is categorised as Web Content Mining, Web structure mining and Web usage mining. Web Mining and its various methods of

implementation are also discussed. Various applications of web content mining, web structure mining and web usage mining are also discussed in this unit. Web mining is the process of extracting useful information from the content, structure, and usage of web data. It involves applying data mining techniques to the web to uncover patterns and insights that can help in various domains, such as business intelligence, search engines, and user experience design. Web mining can be divided into three main categories:

1. **Web Content Mining**: This focuses on extracting information from the content of web pages, such as text, images, videos, or structured data like tables. Techniques like natural language processing (NLP), machine learning, and information retrieval are commonly used.

2. **Web Structure Mining**: It analyzes the structure of the web, such as links between pages or websites. This helps in understanding how web pages are interconnected and can be used in applications like search engine optimization (SEO) or discovering communities within a network. Algorithms like PageRank are examples of web structure mining.

3. **Web Usage Mining**: This focuses on analyzing user behavior by mining web server logs, browser histories, or cookies to track user patterns. It helps in personalizing web experiences, improving website design, and targeting marketing efforts.

Web mining combines techniques from data mining, machine learning, and artificial intelligence to make sense of the vast and dynamic data available on the web. Web Mining algorithm, Page Rank algorithm is explained in this unit with an example. To summarise the unit focuses on the following:

➢ Web Mining
➢ Web Content Mining
➢ Web Structure Mining
➢ Web Usage Mining
➢ Page Rank algorithm

## Review Questions

Q1. What is Web Content Mining? How can web content mining be used for sentiment analysis?

Q2. Explain Page Rank algorithm with example.

Q3. What is recommender system?

Q4. How does the web user mining helps to deliver personalized experience?

Q5. What kind of filtering is used in recommender systems?

# UNIT 11 : TEXT MINING

**Structure**

# 11.0 Introduction

The unit focuses on the Text Mining as an essential process of finding relevant and useful knowledge Text mining, also known as text data mining or text analytics, is the process of extracting useful information and insights from unstructured text data. It involves various techniques from natural language processing (NLP), machine learning, and statistics to analyze and interpret textual data. The goal of text mining is to transform raw text into structured data that can be used for various types of analysis, such as discovering patterns, trends, and relationships within the text. The different techniques of text mining are also discussed in this unit. Episode Rule Discovery for text is explained in the chapter. Hierarchy of categories and their structure is discussed in the chapter.

# 11.1 Objectives

After the end of this unit, you should be able to:

1. Explain Text Mining

2. Explain different techniques in Text Mining

3. Understand Episode Rule Discovery for text

4. Understand Hierarchy of Categories

5. Understand Sentiment Analysis

# 11.2 Text Mining

Text mining, also known as text analytics, is the process of deriving valuable insights and knowledge from unstructured text data. It involves techniques and methods to extract, analyze, and interpret information from textual sources, such as documents, emails, social media posts, customer reviews, and more. Text mining helps uncover patterns, relationships, sentiments, and trends within large volumes of textual data. Some key aspects and techniques used in text mining:

1. Text Pre-processing: Text pre-processing involves cleaning and transforming raw text data to prepare it for analysis. This includes tasks like removing punctuation, converting to lowercase, removing stop words (commonly used words like "and," "the," etc.), stemming or lemmatization (reducing words to their base or root form), and handling special characters or encoding issues.

2. Text Classification: Text classification assigns predefined categories or labels to documents based on their content. It involves training a machine learning model on labelled data and then using it to classify new, unlabeled documents. Common text classification tasks include sentiment analysis, topic classification, spam detection, and document categorization.

3. Named Entity Recognition (NER): NER is a technique used to identify and extract named entities (such as names of people, organizations, locations, dates, etc.) from text. It helps in understanding the key entities mentioned in a document and their relationships.

4. Text Clustering: Text clustering groups similar documents together based on their content or similarity measures. Clustering algorithms, such as k-means or hierarchical clustering, are used to identify patterns or themes within a collection of documents. It helps in organizing and summarizing large document collections and can be used for information retrieval and knowledge discovery.

5. Topic Modelling: Topic modelling is a statistical technique used to discover underlying topics or themes within a collection of documents. Algorithms like Latent Dirichlet Allocation (LDA) or Non-Negative Matrix Factorization (NMF) are commonly used for topic modelling. It helps in understanding the main topics discussed in a set of documents without the need for predefined categories.

6. Sentiment Analysis: Sentiment analysis aims to determine the sentiment or opinion expressed in a piece of text. It involves classifying text as positive, negative, or neutral. Sentiment analysis techniques can be used to analyze customer feedback, social media posts, product reviews, and other forms of text data to understand customer sentiment and opinions.

7. Text Summarization: Text summarization techniques aim to generate concise summaries of long documents or text passages. It can be done using extractive approaches that select important

sentences or phrases from the original text or abstractive approaches that generate new sentences to capture the main ideas.

---

***Check Your Progress 1***

1. What is Text Mining?

2. What is Sentiment Analysis?

3. What is the usage of NER techniques?

4. List major applications of text mining.

5. What is Text Clustering?

---

## 11.3 Text Mining And Unstructured Text

Text mining is a field of study that deals with extracting useful information and knowledge from unstructured text data. Unstructured text refers to textual data that does not have a predefined format or structure, such as emails, social media posts, customer reviews, articles, or any text that is not stored in a structured database. Unstructured text presents unique challenges in terms of its analysis and extraction of valuable insights. Unlike structured data, which is organized in tables or fields, unstructured text lacks a predefined schema, making it difficult to process and analyze using traditional data analysis techniques. Text mining techniques are designed to address these challenges and enable meaningful analysis of unstructured text. Some essential aspects of text mining in relation to unstructured text:

1. Text Pre-processing: Unstructured text data often requires preprocessing before analysis. This involves tasks such as tokenization (splitting text into words or phrases), removing stop words (commonly used words with little semantic value), stemming or lemmatization (reducing words to their base form), handling special characters or punctuation, and encoding or normalizing the text. Pre-processing helps in cleaning and transforming the unstructured text into a structured format suitable for analysis.

2. Information Extraction: Text mining techniques are applied to extract specific information from unstructured text. This can include extracting named entities (such as names, locations, dates), identifying key phrases or topics, detecting sentiment or opinion, extracting relationships between entities, or even identifying events or trends. Techniques like named entity recognition, topic modelling, sentiment analysis, and relation extraction are commonly used for information extraction from unstructured text.

3. Text Classification and Clustering: Text mining involves categorizing or clustering unstructured text into meaningful groups. Text classification assigns predefined categories or labels to text documents based on their content. Clustering algorithms group similar documents together based on their similarity. These techniques help in organizing and structuring unstructured text, enabling better information retrieval and understanding of text patterns and themes.

4. Text Summarization: Unstructured text often contains a vast amount of information, making it challenging to comprehend or extract key insights. Text summarization techniques aim to generate concise summaries of long documents or passages, capturing the main ideas or information. Summarization can be extractive, where important sentences or phrases are selected from the original text, or abstractive, where new sentences are generated to summarize the content.

5. Sentiment Analysis: Sentiment analysis focuses on determining the sentiment or opinion expressed in unstructured text. It involves classifying the text as positive, negative, or neutral to understand the sentiment conveyed. Sentiment analysis is useful in analyzing customer feedback, social media sentiment, product reviews, or public opinion on specific topics. Text mining techniques for unstructured text data leverage machine learning, natural language processing (NLP), and statistical algorithms to process, analyze, and derive insights from textual information. These techniques help overcome the challenges posed by unstructured text and unlock valuable knowledge and information hidden within the vast amount of textual data.

---

*Check Your Progress 2*

Q1. What is Text classification and Text Clustering?

Q2. What is Text Summarization?

Q3. List major steps in text mining.

Q4. What do you understand by unstructured text?

---

## 11.4 Episode Rule Discovery For Text

The concept of "episode rule discovery" is commonly associated with sequential pattern mining, which is a technique used to find interesting patterns or sequences of events in sequential data. While sequential

pattern mining is traditionally applied to transactional data or event sequences, it can also be adapted for text data analysis, where the "episodes" represent patterns of words or phrases occurring in a specific order within textual sequences. Episode rule discovery can be applied to text data in the following ways:

1.  Define the Episodes: In the context of text mining, episodes can be defined as sequences of words or phrases that occur together in a specific order within textual data. For example, an episode might be a specific phrase or a sequence of words that commonly appears in a particular context, such as "customer satisfaction" followed by "improvement" or "product quality" followed by "complaints."

2.  Pre-process the Text Data: Before applying episode rule discovery, the text data needs to be pre-processed. This typically involves steps such as tokenization (splitting text into individual words or phrases), removing stop words, stemming or lemmatization, and handling other text-specific pre-processing tasks.

3.  Identify Sequential Patterns: Sequential patterns refer to regularities or recurring sequences in data, where the order of events or items is crucial. In data mining, sequential pattern mining is the process of finding such patterns in a sequence database, which is useful in various domains like market basket analysis, bioinformatics, web usage mining, and more. A sequential pattern mining algorithm, such as the Prefix Span algorithm or the Apriori algorithm with sequence extension is applied to discover sequential patterns or episodes in the pre-processed text data. These algorithms search for frequent patterns of words or phrases occurring in a specific order across the text documents.

4.  Define Constraints: Specify constraints or interestingness measures to filter and extract meaningful episode rules. These constraints can include minimum support thresholds (minimum occurrence frequency), minimum length or size of episodes, or other domain-specific criteria to focus on relevant patterns.

5.  Evaluate and Interpret Results: Analyze the discovered episode rules to understand their significance and interpret the patterns within the context of the text data. Consider factors like support (frequency of occurrence), confidence (conditional probability), or lift (deviation from expected probability) to assess the importance or interestingness of the discovered patterns.

6.  Apply in Text Mining Tasks: Utilize the discovered episode rules in various text mining tasks, such as information retrieval, sentiment analysis, content recommendation, or knowledge discovery. The patterns can provide insights into recurring patterns of words or phrases that are useful for understanding relationships, identifying trends, or extracting meaningful information from the text data.

Text clustering, also known as document clustering or text categorization, is the process of grouping similar text documents into clusters based on their content or similarity. It is a common technique used in text mining and natural languages processing to organize large collections of documents, identify thematic groups, and facilitate information retrieval and knowledge discovery. The text clustering process includes the following:

1. **Pre-process the Text Data:** Before performing text clustering, the text data needs to be pre-processed. This typically involves steps such as tokenization (splitting text into individual words or phrases), removing stop words, stemming or lemmatization, handling special characters, and transforming the text into a numerical representation suitable for clustering algorithms (e.g., bag-of-words, TF-IDF).

2. **Select a Clustering Algorithm:** Choose an appropriate clustering algorithm that suits the nature of your text data and the objectives of your analysis. Commonly used algorithms for text clustering include k-means clustering, hierarchical clustering, DBSCAN (Density-Based Spatial Clustering of Applications with Noise), and affinity propagation.

3. **Define Similarity Measure:** Select a similarity measure or distance metric to quantify the similarity between pairs of text documents. Commonly used similarity measures for text data include cosine similarity, Euclidean distance, Jaccard similarity, or TF-IDF weighted distance. The choice of similarity measure depends on the characteristics of the text data and the specific requirements of the clustering task.

4. **Apply Clustering Algorithm:** Apply the selected clustering algorithm to the preprocessed text data using the defined similarity measure. The algorithm will group similar documents into clusters based on their content similarity. The number of clusters may be predefined or determined automatically based on the algorithm or evaluation metrics.

5. **Evaluate and Interpret Results:** Evaluate the quality of the generated clusters using appropriate metrics such as silhouette score, cohesion, separation, or purity. Assess the coherence and distinctiveness of the clusters to determine if they align with the expected patterns or themes in the text data. Interpret the resulting clusters to gain insights and understand the underlying patterns or topics present in the data.

6. **Refine and Iteratively Improve:** Iterate and refine the clustering process as needed. This may involve adjusting preprocessing steps, experimenting with different clustering algorithms or

parameter settings, or incorporating domain-specific knowledge to improve the quality and interpretability of the clusters.

7. **Utilize Clusters:** Once the text data is clustered, the resulting clusters can be utilized for various purposes such as information retrieval, topic analysis, recommendation systems, sentiment analysis, or knowledge discovery. The clusters help organize and structure the text data, enabling efficient retrieval of relevant documents or insights from specific thematic groups.

Text clustering is a versatile technique that can be applied to a wide range of text mining tasks. It aids in organizing and understanding large volumes of textual data, enabling efficient analysis and extraction of meaningful information. One popular text clustering algorithm is the k-means clustering algorithm. It is a partition-based algorithm that groups similar data points into a predetermined number of clusters, with the aim of minimizing the within-cluster sum of squares. The algorithm can be adapted for text clustering by representing text documents as numerical vectors using techniques like bag-of-words or TF-IDF. The application of k-means algorithm to text clustering comprises of the following phases-

1. Data Pre-processing: Pre-process the text data by removing stop words, punctuation, and performing stemming or lemmatization. Convert the preprocessed text into a numerical representation, such as a document-term matrix using bag-of-words or TF-IDF.

2. Select the Number of Clusters: Decide on the desired number of clusters, denoted as 'k', based on the specific requirements of the analysis or domain knowledge.

3. Initialize Centroids: Randomly select 'k' initial centroids from the data points. These centroids represent the initial cluster centers.

4. Assign Data Points to Clusters: Calculate the similarity between each data point and the centroids using a distance metric such as Euclidean distance or cosine similarity. Assign each data point to the cluster with the closest centroid.

5. Update Centroids: Recalculate the centroids for each cluster by taking the mean of the feature values (word frequencies or TF-IDF values) of the data points within that cluster.

6. Repeat Steps 4 and 5: Iteratively reassign data points to clusters based on the updated centroids and recalculate the centroids until convergence. Convergence occurs when the centroids no longer change significantly or when a maximum number of iterations is reached.

7. Evaluation: Evaluate the quality of the clustering results using evaluation metrics such as silhouette score, cohesion, or separation. These metrics assess the compactness and separation of the clusters to determine the effectiveness of the clustering algorithm.

8. Interpretation and Analysis: Analyze the resulting clusters to gain insights and interpret the patterns present in the text data. Explore the content of the clusters to understand the themes or topics represented by each cluster.

The k-means algorithm requires the number of clusters ('k') to be predefined. Choosing an appropriate value for 'k' can be challenging and may require experimentation or domain knowledge. Additionally, the algorithm may not always converge to the global optimum and can be sensitive to the initial centroid selection. Other text clustering algorithms include hierarchical clustering (agglomerative or divisive), density-based clustering algorithms like DBSCAN (Density-Based Spatial Clustering of Applications with Noise), or probabilistic models like Latent Dirichlet Allocation (LDA). The choice of algorithm depends on the characteristics of the text data, the desired level of interpretability, and the specific objectives of the clustering task.

---

*Check your progress 2*

1. What is Text clustering?

2. Explain the application of Text Clustering in Text mining.

3. Define episode rule discovery .

4. List major steps in text clustering procedure.

---

## 11.6 Hiearchy Of Categories

In text mining, a hierarchy of categories refers to a hierarchical structure or taxonomy that organizes text documents into multiple levels of categories. It provides a way to organize and classify textual data based on increasingly specific or granular categories. A hierarchy of categories allows for a more organized and systematic approach to text mining, enabling better information retrieval and analysis. An example to illustrate the hierarchy of categories in text mining:

Consider a collection of news articles. The top-level category could be "News," representing the overall domain. Underneath the "News" category, you might have subcategories such as "Sports," "Politics," "Entertainment," and "Technology." Each of these subcategories can further be divided into more specific categories. For instance, under "Sports," you may have subcategories like "Football,"

"Basketball," and "Cricket." The hierarchical structure allows for a nested arrangement of categories, where each category is more specific than its parent category. This hierarchical organization can continue to multiple levels, depending on the complexity and depth required for the analysis.The advantages of a hierarchy of categories in text mining:

1. Organization: The hierarchy provides a structured and organized way to categorize and manage large volumes of text documents. It helps in organizing and navigating the textual data, making it easier to locate specific information.

2. Granularity: The hierarchical structure allows for different levels of categorization, enabling analysis at various levels of specificity. It allows users to drill down into specific subcategories or topics of interest for more focused analysis.

3. Flexibility: The hierarchical taxonomy can be adapted and expanded as needed. New categories can be added, and existing categories can be modified or refined to accommodate evolving requirements or changes in the domain.

4. Information Retrieval: The hierarchy facilitates efficient information retrieval by allowing users to search for documents within specific categories or at different levels of the hierarchy. Users can navigate through the categories to find relevant documents based on their specific information needs.

5. Insights and Analysis: The hierarchy of categories can provide insights into the distribution and relationships between different topics or themes within the text data. It enables the identification of broader patterns and trends across the categories and facilitates comparative analysis between different levels of granularity.

6. Creating a hierarchy of categories for text mining often involves a combination of manual and automated approaches. Domain experts may define the top-level categories and initial subcategories based on their knowledge and understanding of the domain. Automated techniques like text clustering or topic modelling can also be employed to automatically discover patterns or themes in the text data, which can inform the creation of categories at different levels of the hierarchy.

7. Overall, a hierarchy of categories in text mining enhances the organization, analysis, and retrieval of textual data, enabling more efficient and meaningful insights from large document collections.

# 11.7 Sentiment Analysis

Sentiment analysis, also known as opinion mining, is the process of determining the sentiment or emotion expressed in a piece of text. It involves analyzing and classifying textual data to identify whether the sentiment conveyed is positive, negative, or neutral. Sentiment analysis is widely used in various fields, including market research, social media analysis, customer feedback analysis, brand monitoring, and reputation management. The sentiment analysis process includes the following:

1. Data Preparation: Gather the text data that you want to analyze for sentiment. This can include customer reviews, social media posts, survey responses, or any other form of textual data. Data preparation is a crucial step in the data mining process, ensuring that the data is suitable for analysis and can produce accurate and meaningful results. It includes various methods such as:

- **Data Collection:** It is the process of gathering data from various sources such as databases, web scraping, surveys, or sensors. It ensures that the data sources are relevant to the problem being solved, and the data is of sufficient quality and quantity.
- **Data Cleaning:** It is the process of removing or correcting errors, inconsistencies, and inaccuracies in the data. It includes the following:

    i. **Handling Missing Values**: Replace missing values with mean/median, remove records, or use techniques like interpolation.
    ii. **Outlier Detection**: Identify and handle outliers that may skew the analysis.
    iii. **Noise Reduction**: Smooth noisy data using techniques like binning or clustering.
    iv. **Correcting Data Types**: Ensure all data fields are of the correct type (e.g., numeric, categorical).

- **Data Integration:** It is the process of combining data from different sources into a unified dataset.
- **Data Transformation:** It is the process of modifying the data to make it more suitable for mining tasks. The essential techniques are:

    a. **Normalization/Standardization**: Adjusting the scale of data to ensure uniformity.
    b. **Aggregation**: Summarizing data (e.g., weekly sales totals).
    c. **Dimensionality Reduction**: Reducing the number of variables using techniques like PCA (Principal Component Analysis).
    d. **Feature Selection/Extraction**: Identifying and selecting the most relevant features for analysis.

206

- **Data Reduction:** It is the process of reducing the volume of data while maintaining its integrity. Various techniques include:

  a. **Sampling**: Selecting a representative subset of the data.
  b. **Clustering**: Grouping similar data points to reduce complexity.
  c. **Dimensionality Reduction**: Reducing the number of attributes (features) in the dataset.

- **Data Discretization:** It is the process of converting continuous data into discrete buckets or intervals. It helps in simplifying the data and making it more suitable for certain algorithms, particularly in classification and association rule mining.

2. Text Pre-processing: Pre-process the text data to clean and transform it into a suitable format for sentiment analysis. Steps may include removing punctuation, converting to lowercase, removing stop words, handling special characters or emojis, and normalizing the text. Various approaches to sentiment analysis includes the following:

I. Lexicon-Based Approaches: Lexicon-based approaches use sentiment lexicons or dictionaries that contain pre-defined sentiment scores for words or phrases. Each word is assigned a polarity (positive, negative, or neutral) based on its sentiment score. The sentiment of the text is calculated by aggregating the sentiment scores of individual words in the text. Some popular sentiment lexicons include AFINN, SentiWordNet, and VADER.

II. Machine Learning Approaches: Machine learning algorithms can be used for sentiment analysis, where a model is trained on labeled data (text with known sentiment labels) to classify the sentiment of new, unlabeled text. Common machine learning algorithms used for sentiment analysis include Naive Bayes, Support Vector Machines (SVM), and Recurrent Neural Networks (RNN), particularly Long Short-Term Memory (LSTM) models.

III. Sentiment Classification: Apply the chosen sentiment analysis approach to classify the sentiment of the text. The output can be binary (positive/negative) or multi-class (positive/negative/neutral), depending on the requirements of the analysis. The sentiment classification can be performed at the document level (e.g., classifying an entire customer review) or at the sentence or aspect level (identifying sentiment towards specific aspects or entities mentioned in the text).

3. Evaluation: Evaluate the performance of the sentiment analysis model using appropriate evaluation metrics such as accuracy, precision, recall, or F1 score. This step is crucial to assess the reliability and effectiveness of the sentiment analysis approach and fine-tune the model if necessary.

4. Interpretation and Application: Analyze and interpret the sentiment analysis results to gain insights from the text data. Understand the sentiment distribution, identify trends, or detect sentiment patterns related to specific topics, products, or entities. The insights derived from sentiment analysis can be utilized for decision-making, brand management, customer service improvement, or other business applications.

---

*Check Your Progress 3*

1. What is Sentiment Analysis?

2. What is Sentiment Classification?

3. List the advantages of text mining.

4. Define granularity.

5. What is the purpose of pre-processing required to perform sentiment analysis?

6. Define sentiment lexicon.

---

## 11.9 Summary

The unit discussed the process of Text Mining. Various techniques in Text Mining are also explained in this unit. Text mining plays a crucial role in various domains by enabling the extraction of meaningful information from vast amounts of unstructured text data. Its ability to transform text into actionable insights has significant implications for business, research, healthcare, social sciences, and more. The concept of Episode rule discovery for text is explained in this unit. Hierarchy of concepts and Sentiment Analysis is also discussed in this unit.

To summarise the unit focuses on the following:

- Different techniques in Text Mining

- Episode Rule Discovery for text

- Hierarchy of Categories

- Sentiment Analysis

**Terminal Questions**

1. What is Text Mining? How does Text Mining differ from Data Mining?

2. Discuss some techniques used in Text Mining?

3. What are some common techniques used in Text Mining?

4. What are the ethical considerations in Text Mining?

5. Discuss the various Data preparation methods in Text mining process.

# UNIT 12 SPATIAL MINING

**Structure**

## 12.0 Introduction

The unit focuses on the concept of Spatial Mining. The unit discusses various essential methods and algorithms in Spatial Mining. Spatial data mining means extraction of knowledge, spatial relationships, or other interesting patterns stored in spatial databases or other databases and data repositories. The unit explains spatial data mining as the process of discovering interesting and previously unknown patterns, relationships, and trends from large spatial datasets. These datasets include information that has a geographical or spatial component, such as maps, satellite images, GPS data, and other forms of location-based data. Spatial data mining extends traditional data mining techniques by incorporating spatial attributes, making it possible to analyze data with respect to its location and spatial context.

## 12.1 Objectives

After the end of this unit, you should be able to:

1.  Explain the concept of Spatial Mining

2.  Explain various methods in Spatial Mining

3.  Understand Spatial Clustering

4.  Understand Spatial Clustering methods

## 12.2 Spatial Mining

Spatial mining, also known as spatial data mining or geographic data mining, is the process of discovering interesting and useful patterns, relationships, or knowledge from spatial or geographic data. It involves analyzing data that has an inherent spatial component, such as geographic coordinates, regions, distances, or spatial relationships. Spatial mining techniques combine principles from data mining, spatial analysis, and geographic information systems (GIS) to extract valuable insights from spatial data. Here are some key concepts and techniques related to spatial mining:

1. Spatial Data Representation: Spatial data can be represented in various formats, including point data (individual coordinates), line data (e.g., road networks), polygon data (e.g., administrative boundaries), or raster data (grids or images). Effective representation of spatial data is crucial for spatial mining.

2. Spatial Patterns and Relationships: Spatial mining aims to identify patterns and relationships within spatial data. This can include spatial clustering (grouping of similar spatial objects), spatial autocorrelation (tendency of nearby objects to have similar attributes), spatial outliers (unusual or exceptional spatial objects), spatial association (correlations between spatial objects), or spatial dependence (influence of one spatial object on others).

3. Spatial Data Preprocessing: Preprocessing spatial data involves tasks such as data cleaning, data integration (combining multiple sources of spatial data), data transformation (projection, scaling), or feature extraction (extracting relevant spatial attributes). Preprocessing is essential for preparing spatial data for analysis.

4. Spatial Data Mining Techniques: Spatial data mining techniques adapt traditional data mining algorithms and methods to handle spatial data. These techniques can include clustering algorithms (e.g., DBSCAN, OPTICS), classification algorithms (e.g., decision trees, support vector machines), association rule mining, spatial regression, spatial interpolation, or spatial outlier detection.

5. Spatial Querying and Retrieval: Spatial mining often involves querying and retrieving specific spatial data based on location or spatial attributes. Spatial query languages, such as SQL with spatial extensions or specialized query languages like Spatial Query Language (SQL), allow users to extract relevant spatial information from databases or GIS systems.

6. Geospatial Visualization: Visualizing spatial data and analysis results is essential for interpreting and communicating spatial patterns and relationships. Geospatial visualization techniques include maps,

spatial heatmaps, choropleth maps, scatter plots, or 3D visualization methods, which help users understand and explore the spatial mining results.

Spatial mining encompasses various topics and areas of research. Here are some key topics in spatial mining:

1. Spatial Data Clustering: This topic focuses on clustering algorithms specifically designed for spatial data, as discussed earlier. It includes techniques such as DBSCAN, OPTICS, spatial agglomerative clustering, and others, which aim to identify spatially coherent clusters within spatial datasets.

2. Spatial Association Analysis: Spatial association analysis examines the relationships and dependencies between spatial objects or attributes. It involves identifying spatial patterns, co-occurrence of events or phenomena, and quantifying the degree of association between spatial variables. Techniques like spatial autocorrelation, spatial dependence models, and association rule mining can be used in this context.

3. Spatial Trajectory Mining: This topic deals with analyzing and mining patterns from moving object trajectories. It involves techniques for trajectory segmentation, clustering, outlier detection, trajectory pattern discovery, and prediction. Spatial trajectory mining finds applications in transportation, mobile devices, animal tracking, and other domains where movement data is available.

4. Spatial Data Classification: Spatial data classification focuses on assigning spatial objects to predefined classes or categories based on their spatial and attribute characteristics. It includes supervised learning methods, such as decision trees, support vector machines, and random forests, which utilize spatial features and attributes for classification tasks.

5. Spatial Outlier Detection: Spatial outlier detection aims to identify spatial objects or regions that deviate significantly from the expected spatial distribution or patterns. It involves techniques such as distance-based outliers, local outliers, global outliers, and spatial anomaly detection. Spatial outlier detection is important in various domains, including environmental monitoring, urban planning, and crime analysis.

6. Spatial Data Visualization: Visualization plays a crucial role in spatial mining to explore and present spatial patterns and relationships effectively. It includes techniques for mapping spatial data, creating thematic maps, spatial heatmaps, choropleth maps, 3D visualization, and interactive visual analytics tools.

7. Spatiotemporal Data Mining: Spatiotemporal data mining focuses on analyzing and discovering patterns in data that have both spatial and temporal dimensions. It involves techniques for spatiotemporal clustering, trajectory analysis, change detection, event detection, and spatiotemporal

pattern mining. This area finds applications in areas like transportation, environmental monitoring, public health, and social media analysis.

8. Geospatial Data Integration: Geospatial data integration deals with combining and integrating multiple sources of spatial data from different formats or coordinate systems. It involves spatial data fusion, data harmonization, spatial data matching, and spatial data integration frameworks to create unified and consistent spatial datasets for mining and analysis.

Applications of spatial mining are diverse and can be found in various domains, including urban planning, transportation, environmental monitoring, epidemiology, market analysis, crime analysis, and natural resource management. Spatial mining techniques enable the discovery of valuable insights, support decision-making processes, and facilitate understanding of spatial relationships and phenomena. It's important to note that spatial mining requires specialized knowledge and tools to handle the unique characteristics and challenges of spatial data. Geographic Information Systems (GIS) software and spatial databases often provide built-in functionalities for spatial data mining and analysis.

## 12.3. Spatial Mining Clustering

Spatial mining clustering refers to the application of clustering algorithms specifically designed for spatial or geographic data. These algorithms aim to identify groups or clusters of spatial objects that exhibit similar characteristics or spatial proximity. Spatial clustering is useful for discovering patterns, hotspots, or regions of interest within spatial data. Some spatial mining algorithms are:

1. DBSCAN (Density-Based Spatial Clustering of Applications with Noise): DBSCAN is a popular density-based clustering algorithm. It groups spatial objects based on their density and identifies dense regions separated by sparse areas. It is capable of finding clusters of arbitrary shape and can handle noisy data effectively.

2. OPTICS (Ordering Points To Identify the Clustering Structure): OPTICS is an extension of DBSCAN that provides a hierarchical view of the clustering structure. It orders points based on their reachability distancegr, enabling the identification of clusters at different density levels.

3. K-means: K-means is a well-known partition-based clustering algorithm. In spatial mining, it can be adapted by considering the spatial coordinates of objects along with their attribute values. The resulting clusters can represent spatially coherent groups.

4. Spatial Agglomerative Clustering: This hierarchical clustering approach starts with each spatial object as a separate cluster and iteratively merges clusters based on a distance or similarity measure. It creates a hierarchy of clusters that can be explored at different levels of granularity.

5. ST-DBSCAN (Spatial-Temporal DBSCAN): ST-DBSCAN is an extension of DBSCAN that incorporates both spatial and temporal dimensions. It can identify clusters that exist over both space and time, making it suitable for analyzing spatiotemporal data.

6. Spatial Scan Statistics: Spatial scan statistics identify statistically significant spatial clusters by comparing the observed spatial distribution of events with the expected distribution under a null hypothesis. These methods are widely used in epidemiology, crime analysis, and other domains where identifying clusters of events is important.

These are just a few examples of spatial clustering techniques used in spatial mining. The choice of clustering algorithm depends on the specific characteristics of the spatial data, the desired properties of the clusters (e.g., shape, size, density), and the goals of the analysis. It's also worth noting that spatial clustering algorithms often take into account spatial distances or neighbourhood relationships to capture spatial proximity and structure within the data. Spatial mining algorithms refer to techniques specifically designed to extract valuable patterns, relationships, or knowledge from spatial or geographic data.

---

*Check Your Progress*

1. What is spatial mining?

2. What is spatial mining clustering?

3. Explain some methods in spatial mining.

4. List some essential methods in Spatial Mining.

---

## 12.5 Summary

The unit explained the concept of Spatial Mining. "Spatial mining" refers to the process of extracting and analyzing spatial data to identify patterns, relationships, or insights. Spatial mining, also known as spatial data mining, is the process of discovering interesting patterns and relationships in spatial data. Spatial data refers to data that is related to locations in a geographic space, such as maps, satellite images, GPS data, and location-based services. The unique aspect of spatial mining is that it not only considers the traditional attributes of data but also incorporates the spatial relationships and structures inherent in the data. This type of data is often related to geographic locations, and spatial mining techniques are commonly used in fields such as geographic information systems (GIS), urban planning, environmental monitoring, and remote sensing. The unit explained the various essential methods and algorithms in Spatial Mining.

# Terminal Questions

1. What are the primary techniques used in spatial data mining to handle the spatial autocorrelation inherent in spatial datasets, and how do these techniques ensure the validity and accuracy of the results?

2. Describe the process of spatial clustering and its applications in real-world scenarios. How do algorithms like DBSCAN or K-means need to be adapted when applied to spatial data?

3. Explain the concept of spatial outlier detection and its importance in spatial data mining. What methods are commonly used to identify spatial outliers, and what are the practical applications of detecting such outliers?

4. Discuss the role of spatial data mining in disaster management and response. How can spatial mining techniques be used to enhance the preparedness and response to natural disasters such as earthquakes, floods, or hurricanes?

**Note**