**Dimensionality Reduction:** Unsupervised: Introduction, Subset Selection, PCA (Principal Component Analysis) – Technique, Examples as Numerical. Mining Different types of data: Mining the World Wide Web - Page Rank Algorithm, Text mining, Mining Time Series Data, Ensemble methods-Increasing the Accuracy.

**Textbooks:**

1. Jiawei Han and Micheline Kamber: Data Mining Concepts and Techniques,Elsevier, 2nd Edition, 2009.
2. Stephen Marsland, "Machine Learning - An Algorithmic Perspective", Second Edition, CRC Press - Taylor and Francis Group, 2015.
3. Ethem Alpaydin, "Introduction to Machine Learning", Second Edition, MITPress, Prentice Hall of India (PHI) Learning Pvt. Ltd. 2010.
4. Xindong Wu and Vipin Kumar: The top ten Algorithms in Data Mining, Chapman and Hall/CRC press.
5. Pang-Ning Tan, Michael Steinbach and Vipin Kumar, "Introduction to Data Mining", Pearson Education, 2007.
6. DISCOVERING KNOWLEDGE IN DATA, An Introduction to Data Mining, Second Edition, Daniel T. Larose ,Chantal D. Larose.

**Reference Books:**

1. K.P. Soman, ShyamDiwakar and V. Aja, "Insight into Data Mining Theory and Practice", Eastern Economy Edition, Prentice Hall of India, 2006.
2. G. K. Gupta, "Introduction to Data Mining with Case Studies", Eastern Economy Edition, Prentice Hall of India, 2006.
3. Christopher Bishop, "Pattern Recognition and Machine Learning", CBS Publishers & Distributors, 2010.
4. Mehryar Mohri, Afshin R, Ameet Talwalkar, "Foundations of Machine Learning", MIT Press, 2012.

1. **Explain the working of Google page rank algorithm, formula and damping factor.**

The Google PageRank algorithm is a key component of Google's search engine algorithm that determines the importance or relevance of web pages. It was developed by Google's co-founders Larry Page and Sergey Brin and has since become one of the most influential ranking algorithms on the web. While the specific details of Google's algorithm are proprietary and not publicly disclosed, I can provide a general overview of how PageRank works.

The basic idea behind PageRank is to assign a numerical value, known as a PageRank score, to each web page on the internet. This score represents the relative importance of a page in the overall link structure of the web. The underlying assumption is that a page is considered more important if it is linked to by other important pages.

The algorithm works through an iterative process that starts by assigning an initial PageRank score to each web page. This initial score can be evenly distributed among all pages or based on certain predefined criteria. Let's call this initial score PR(A) for a page A.

The algorithm then goes through a series of iterations, where it recalculates the PageRank score for each page based on the scores of the pages that link to it. In other words, the score of a page is influenced by the scores of the pages that point to it. The formula for calculating the updated PageRank score for a page A is as follows:

$$PR(A) = (1 - d) + d * (PR(T1)/C(T1) + PR(T2)/C(T2) + ... + PR(Tn)/C(Tn))$$

In this formula:

- PR(A) is the new PageRank score for page A.
- PR(T1), PR(T2), ..., PR(Tn) are the PageRank scores of the pages T1, T2, ..., Tn that link to page A.
- C(T1), C(T2), ..., C(Tn) are the total number of outbound links from the pages T1, T2, ..., Tn.
- d is the damping factor, a value between 0 and 1 that represents the probability that a user will continue clicking on links rather than abandoning the web. Google typically uses a damping factor of around 0.85.

The damping factor is an important component of PageRank as it ensures that the algorithm does not get stuck in an infinite loop. It models the behavior of a random web surfer who, at each page, either clicks on a link with probability d or randomly jumps to any page on the web with probability (1 - d).

During the iterative process, the PageRank scores are recalculated until they converge,

meaning that the scores stabilize and do not change significantly between iterations. Once the convergence is reached, the final PageRank scores represent the relative importance of each page in the web graph.

2. **Explain the following:**
   i. **Web mining**
   ii. **Usage mining**
   iii. **Link mining**
   iv. **Statistical relational learning**
   v. **Web content mining.**

i. **Web mining:**

Web mining refers to the process of extracting valuable information and knowledge from web data. It involves the use of data mining techniques to discover patterns, trends, and relationships within web content, structure, and usage. Web mining can be categorized into three main types: web content mining, web structure mining, and web usage mining. It aims to uncover hidden insights from web resources, such as web pages, hyperlinks, user behavior, and server logs.

ii. **Usage mining:**

Usage mining, also known as web usage mining or clickstream analysis, is a subset of web mining that focuses on analyzing user behavior and interactions on websites. It involves the collection and analysis of web usage data, such as clickstreams, navigation paths, and session logs. Usage mining aims to understand user preferences, identify browsing patterns, and improve website usability. The extracted knowledge can be used for personalization, recommendation systems, targeted advertising, and website optimization.

iii. **Link mining:**

Link mining, also referred to as hyperlink analysis or web structure mining, involves the analysis of the relationships between web pages through hyperlinks. It focuses on studying the link structure and properties of the World Wide Web. Link mining techniques can be used to discover important web pages, determine the relevance and authority of web pages, detect communities or clusters of related pages, and identify spam or malicious websites. Link analysis

algorithms, such as PageRank and HITS (Hyperlink-Induced Topic Search), are commonly used in link mining.

iv. **Statistical relational learning:**

Statistical relational learning (SRL) is a field that combines statistical modeling and machine learning techniques with formal logic and relational databases. It aims to handle data that has both relational and statistical aspects, such as data with complex structures, multiple interrelated tables, or knowledge represented in graphs. SRL algorithms can reason about relational data, make predictions, and discover patterns while incorporating statistical uncertainty. It finds applications in various domains, including social network analysis, knowledge graph construction, recommender systems, and bioinformatics.

v. **Web content mining:**

Web content mining focuses on extracting useful information and knowledge from the actual content of web pages. It involves analyzing the textual, visual, and multimedia data present on web pages. Web content mining techniques can be used to extract structured information, such as extracting product names and prices from e-commerce websites, sentiment analysis of customer reviews, text categorization, information extraction from news articles, and image or video analysis on the web. It often involves natural language processing, image processing, and machine learning techniques to process and understand web content.

3. **Illustrate how using ensemble methods can lead to increased accuracy.**

Ensemble methods refer to the technique of combining multiple individual models to create a stronger and more accurate predictive model. By leveraging the collective wisdom of diverse models, ensemble methods can improve accuracy and generalization capabilities compared to using a single model alone. Here's an illustration of how ensemble methods can lead to increased accuracy:

1. Diverse Model Creation: To create an ensemble, multiple individual models with different characteristics are trained. These models can vary in terms of algorithms, hyperparameters, training data subsets, or initializations. For example, you could use decision trees, support vector machines, neural networks, or any other suitable models.

2. Independent Training: Each individual model in the ensemble is trained

independently using a subset of the training data. By training on different subsets, the models can capture different patterns and make diverse predictions.

3. Combining Predictions: Once the individual models are trained, their predictions are combined to produce the final ensemble prediction. The combination can be achieved through various techniques such as voting, averaging, or weighted averaging. The idea is to aggregate the predictions of different models to create a more robust and accurate prediction.

4. Reducing Bias and Variance: Ensemble methods are effective at reducing both bias and variance. Bias refers to the error caused by the models' assumptions not capturing the true patterns in the data, while variance refers to the error caused by the models being sensitive to fluctuations in the training data. By combining diverse models, ensemble methods can mitigate these issues. Individual models with high bias might make different types of errors, but when combined, their errors can cancel each other out. Similarly, models with high variance can be smoothed out by aggregating their predictions.

5. Improved Generalization: Ensemble methods often lead to improved generalization, allowing the model to perform well on unseen data. The ensemble can capture a broader range of patterns and provide a more robust representation of the underlying data distribution. It can identify complex relationships, handle noisy data, and overcome overfitting issues that individual models may suffer from.

6. Enhanced Stability and Robustness: Ensemble methods are more stable and less prone to random fluctuations in the data compared to single models. As individual models are combined, the ensemble prediction tends to be more reliable and consistent. It can handle outliers, errors, or missing values more gracefully.

7. Increased Accuracy: The ultimate outcome of using ensemble methods is typically an improvement in prediction accuracy. By combining the strengths of multiple models and mitigating their weaknesses, ensemble methods can achieve higher accuracy than any individual model. This increased accuracy can have significant practical implications, especially in domains where accurate predictions are crucial, such as finance, healthcare, or image recognition.

**4. Explain Challenges of Genomic Data Mining.**

Genomic data mining refers to the process of extracting useful information and insights from large-scale genomic datasets. While genomic data has the potential to revolutionize medicine and biology, it also presents several challenges that need to be addressed. Here are some of the key challenges in genomic data mining:

I.   Data Volume and Complexity: Genomic data is vast and complex. Whole-genome sequencing generates terabytes of data per individual, and large-scale studies involve thousands or even millions of individuals. Managing, storing, and processing such large volumes of data requires advanced computational infrastructure and algorithms.

II.  Data Quality and Variability: Genomic data is prone to errors and variability. Sequencing technologies can introduce errors, and biological factors such as DNA damage or contamination can affect data quality. Additionally, genetic variations among individuals, including single nucleotide polymorphisms (SNPs), structural variants, and copy number variations, contribute to the variability of the data.

III. Privacy and Ethical Concerns: Genomic data contains highly sensitive information about individuals, including their genetic predispositions to diseases, ancestry, and potentially identifiable information. Preserving privacy while enabling data sharing and analysis poses significant challenges. Ensuring proper data anonymization, consent management, and secure data storage and transmission are crucial in genomic data mining.

IV.  Integration and Interpretation: Genomic data mining often requires integration with other types of biological and clinical data to gain meaningful insights. Integrating genomic data with phenotypic, environmental, and lifestyle data is essential for understanding the complex interplay between genes and traits. However, integrating diverse data sources with different formats and standards is a significant challenge.

V.   Computational Methods and Tools: Developing efficient algorithms and computational tools for analyzing genomic data is a constant challenge. The complexity and scale of genomic data require advanced statistical and machine learning techniques, as well as scalable and distributed computing frameworks. Furthermore, interpreting the results of genomic data analysis in a biologically meaningful way remains a challenge due to the vastness of our current knowledge of gene function and interactions.

VI.  Data Sharing and Collaboration: Genomic data mining benefits from collaboration

and data sharing across research institutions and countries. However, challenges such as data access, data harmonization, and establishing data sharing agreements hinder widespread collaboration. Overcoming legal, ethical, and technical barriers to data sharing is essential for maximizing the potential of genomic data mining.

VII.   Translation to Clinical Applications: Translating genomic discoveries into clinical applications is a major challenge. Validating the predictive power of genomic markers, identifying clinically actionable variants, and demonstrating the clinical utility of genomic-based interventions require rigorous study designs, large and diverse patient cohorts, and robust statistical analyses.

**5. Explain the following:**
**i)      Mining time series data.**
**ii)     Text mining.**

**i)      Mining time series data:**

- Mining time series data refers to the process of extracting valuable insights, patterns, and knowledge from a collection of data points that are ordered in time. Time series data is characterized by its sequential nature, where each data point is associated with a specific timestamp or temporal information. Examples of time series data include stock prices, weather data, sensor readings, and financial transactions.

- The mining of time series data involves various techniques and algorithms to analyze and discover meaningful patterns and trends within the data. This can be done for several purposes, such as forecasting future values, detecting anomalies or outliers, identifying recurring patterns, and understanding the underlying relationships or dependencies between different variables.

- Common methods used in mining time series data include statistical techniques like time series decomposition, autoregressive integrated moving average (ARIMA) models, and exponential smoothing. Additionally, machine learning approaches such as recurrent neural networks (RNNs), Hidden Markov Models (HMMs), and Long Short-Term Memory (LSTM) networks are also applied for more complex time series analysis tasks.

**iii)    Text mining:**

- Text mining, also known as text analytics or text data mining, involves the process of extracting valuable insights and knowledge from unstructured textual data. Unstructured text data refers to any form of written or recorded information that does

not have a predefined format or structure, such as emails, social media posts, customer reviews, news articles, and documents.

- Text mining utilizes various natural language processing (NLP) techniques and computational algorithms to analyze, interpret, and extract meaningful information from text data. It involves several tasks, including text classification, sentiment analysis, topic modeling, named entity recognition, and information extraction.

- Text mining techniques often involve preprocessing steps such as tokenization (splitting text into individual words or tokens), stemming (reducing words to their base or root form), and removing stop words (common words like "the," "is," etc.). After preprocessing, statistical and machine learning algorithms can be applied to analyze the text data and uncover patterns, trends, or relationships.

- Applications of text mining are diverse and can be found in various fields, including customer feedback analysis, market research, social media monitoring, content recommendation systems, fraud detection, and information retrieval. By extracting valuable insights from unstructured text data, organizations can make informed decisions, improve their products or services, and gain a better understanding of customer opinions and preferences.

**6. Define boosting. Represent it schematically. Write the AdaBoost algorithm pseudocode for two class classification problem. Comment on algorithm performance.**

Boosting is a machine learning technique that combines multiple weak learners to create a strong learner. It is an iterative process where each weak learner focuses on the examples that were misclassified by the previous learners. The final prediction is made by combining the predictions of all the weak learners, typically through a weighted majority voting scheme.

Schematic representation of the AdaBoost algorithm:

1. Initialize the weights of the training
   examples: $w\_1 = 1/N, w\_2 = 1/N, ..., w\_N = 1/N$

2. for t = 1 to T:  // Number of weak learners

3. Train a weak learner on the training data with weights $w\_1, w\_2, ..., w\_N$

4. Calculate the error of the weak learner:
   $epsilon\_t = sum(w\_i * I(y\_i != h\_t(x\_i)))$, for i = 1 to N

5. Calculate the weight of the weak learner:
   $alpha\_t = 0.5 * ln((1 - epsilon\_t) / epsilon\_t)$

6. Update the weights of the training examples:

   w_i = w_i * exp(-alpha_t * y_i * h_t(x_i)), for i = 1 to N

7. Normalize the weights:

   w_i = w_i / sum(w_i), for i = 1 to N

8. end for

In the pseudocode above, N represents the number of training examples, T represents the number of weak learners, h_t(x_i) represents the prediction of the t-th weak learner for the i-th training example, y_i represents the true label of the i-th training example, and I(condition) is an indicator function that returns 1 if the condition is true and 0 otherwise.

The AdaBoost algorithm is known for its ability to improve the performance of weak learners by focusing on misclassified examples in each iteration. The algorithm assigns higher weights to misclassified examples, forcing subsequent weak learners to pay more attention to them. By combining the predictions of multiple weak learners with different weights, AdaBoost creates a strong classifier that can handle complex decision boundaries.

The performance of the AdaBoost algorithm is generally quite good. It has been shown to have good generalization capabilities, often outperforming individual weak learners and other ensemble methods. However, AdaBoost can be sensitive to noisy data and outliers, which can affect its performance. Additionally, if the weak learners are too complex or overfit the data, AdaBoost may suffer from overfitting as well. Proper parameter tuning and selection of appropriate weak learners are important for achieving optimal performance with AdaBoost.

**7. Explain Text Summarization, Document Extraction and Information retrieval in Text Mining.**

**Text Summarization:**

Text summarization is the process of condensing a piece of text into a shorter version while preserving its key information and meaning. It aims to provide a concise and coherent summary of the main ideas and important details contained in a larger text, such as an article, document, or web page. Text summarization can be done in two main ways:

1. Extractive Summarization: In extractive summarization, the summary is created by selecting and extracting the most relevant sentences or phrases from the original text. The extracted sentences are then combined to form a summary. This approach relies

on identifying the most informative and important parts of the text based on various criteria such as sentence relevance, importance, and redundancy.

2. Abstractive Summarization: Abstractive summarization involves generating a summary that may contain new sentences or phrases that are not present in the original text. It goes beyond simple sentence extraction and aims to create a coherent and concise summary by understanding the meaning and context of the text. This approach utilizes natural language processing (NLP) techniques, such as semantic analysis, language generation, and paraphrasing, to generate a summary that captures the essence of the original text.

**Document Extraction:**

Document extraction, also known as information extraction, is the process of automatically identifying and extracting specific pieces of information from unstructured or semi-structured documents. The goal is to transform the textual content of documents into structured data that can be analyzed, organized, and utilized in various applications. Some common types of information that can be extracted from documents include names, dates, locations, organizations, and numerical data.

**Document extraction involves several steps, including:**

1. Text Parsing: Parsing the document to identify its structure, such as sections, paragraphs, and headings, and breaking it down into manageable units for analysis.

2. Entity Recognition: Identifying and extracting named entities, such as persons, organizations, locations, and dates, from the text. This is typically done using NLP techniques, such as named entity recognition (NER).

3. Relationship Extraction: Discovering and extracting relationships between entities mentioned in the document. For example, identifying that a person works for a specific organization or that a company is located in a particular city.

4. Event Extraction: Extracting information about events mentioned in the document, such as conferences, meetings, or product launches.

**Information Retrieval in Text Mining:**

Information retrieval (IR) in text mining refers to the process of searching and retrieving relevant information from a large collection of text documents or sources. It involves techniques and algorithms to efficiently retrieve documents that are most likely to be relevant to a user's query or information need.

The key steps in information retrieval include:

1. Indexing: Preprocessing and indexing the documents to create an inverted index, which is a data structure that maps each term in the document collection to the documents that contain it. This allows for fast lookup and retrieval of documents based on search queries.

2. Query Processing: Analyzing and processing user queries to identify the most relevant terms and constructs. This involves techniques such as query expansion, term weighting, and relevance ranking.

3. Ranking and Scoring: Assigning a relevance score to each document based on its similarity to the query. Various ranking algorithms, such as TF-IDF (Term Frequency-Inverse Document Frequency) and BM25 (Best Match 25), are commonly used to determine the relevance of documents.

4. Retrieval and Presentation: Retrieving and presenting the most relevant documents to the user based on their query. This can involve displaying the documents in a ranked list, highlighting relevant sections, or providing snippets that contain the queried terms.

Information retrieval in text mining is widely used in search engines, recommendation systems, content analysis, and other applications that involve searching and retrieving information from large text collections.

**8. Explain Subset Selection, PCA (Principal Component Analysis) – Techniques**

Subset Selection: Subset selection is a technique used in machine learning and statistics to identify a subset of features or variables that are most relevant for a particular task. In many real-world problems, the available data may contain numerous features, but not all of them may be useful or necessary for the task at hand. Subset selection aims to identify the most informative subset of features that can effectively represent the underlying patterns in the data.

There are two main types of subset selection techniques: forward selection and backward elimination.

1. Forward Selection: Forward selection starts with an empty set and iteratively adds features to the subset one at a time. At each iteration, the algorithm evaluates the performance of the model using a criterion such as accuracy or error rate. The feature that improves the model's performance the most is selected and added to the subset. This process continues until a stopping criterion is met, such as reaching a desired number of features or a decrease in performance.

2. Backward Elimination: Backward elimination starts with all the features included in

the subset and removes one feature at a time in each iteration. Similar to forward selection, a criterion is used to evaluate the model's performance after removing each feature. The feature whose removal causes the least deterioration in performance is eliminated. The process continues until a stopping criterion is met.

Subset selection techniques are useful for feature engineering and dimensionality reduction. By selecting a subset of features, we can simplify the model, reduce computational complexity, and potentially improve the interpretability and generalization ability of the model.

Principal Component Analysis (PCA): PCA is a widely used technique for dimensionality reduction and data compression. It aims to transform a high-dimensional dataset into a lower-dimensional space while preserving the most important patterns and relationships in the data.

The key idea behind PCA is to find a set of orthogonal axes, called principal components, along which the data varies the most. These principal components are ranked in order of the amount of variance they explain in the data. The first principal component explains the largest amount of variance, the second principal component explains the second largest amount of variance, and so on.

PCA works by constructing these principal components as linear combinations of the original variables. The first principal component is a weighted sum of the original variables, where the weights are chosen to maximize the variance. The subsequent principal components are also linear combinations of the original variables but with the constraint that they are orthogonal to the previous components.

The PCA algorithm involves the following steps:
1. Standardize the data by subtracting the mean and scaling by the standard deviation.
2. Compute the covariance matrix or correlation matrix of the standardized data.
3. Perform an eigendecomposition of the covariance/correlation matrix to obtain the eigenvectors and eigenvalues.
4. Sort the eigenvectors in descending order of their corresponding eigenvalues.
5. Select the top-k eigenvectors based on the amount of variance they explain or choose a desired number of dimensions for the lower-dimensional space.
6. Transform the original data onto the new lower-dimensional space spanned by the selected eigenvectors.

PCA is useful for various tasks, including data visualization, noise reduction, feature extraction, and speeding up machine learning algorithms by reducing the input dimensions. By retaining the most important information while reducing dimensionality, PCA can help in

gaining insights from complex datasets and improving computational efficiency.

### 9. Explain Feature Selection and Feature Extraction in detail.

Feature selection and feature extraction are two important techniques used in machine learning and data analysis to improve the performance and efficiency of models. Both techniques aim to reduce the dimensionality of the data by selecting or creating a subset of relevant features. However, they differ in their approach and the methods used. Let's discuss each technique in detail.

1. **Feature Selection:**

    Feature selection involves selecting a subset of relevant features from the original set of features in the dataset. The goal is to retain the most informative and discriminative features while discarding irrelevant or redundant ones. By eliminating irrelevant features, feature selection helps in reducing noise and improving model accuracy, interpretability, and efficiency. There are three common types of feature selection methods:

a. Filter Methods: Filter methods evaluate the relevance of features based on their statistical properties, such as correlation, mutual information, or significance tests. They rank features using metrics like information gain, chi-square, or correlation coefficient, and select the top-ranked features. These methods are computationally efficient but do not consider the interaction with the learning algorithm.

b. Wrapper Methods: Wrapper methods assess the performance of a learning algorithm using different feature subsets. They use a specific machine learning algorithm to evaluate subsets of features and select the one that produces the best performance. Wrapper methods are computationally expensive since they involve training and evaluating multiple models.

c. Embedded Methods: Embedded methods incorporate feature selection as part of the learning algorithm itself. They aim to select features during the model training process, considering their relevance to the predictive task. Embedded methods are efficient since feature selection is performed within the model training loop. Examples of embedded methods include LASSO (Least Absolute Shrinkage and Selection Operator) and regularization techniques like ridge regression and elastic net.

2. **Feature Extraction:**

    Feature extraction involves creating new features by transforming the original set of features into a lower-dimensional representation. The goal is to capture the essential information contained in the original features while reducing the dimensionality.

Feature extraction techniques are often used when the original features are high-dimensional or when domain knowledge suggests that a different representation may be more meaningful. Here are some common feature extraction methods:

a. Principal Component Analysis (PCA): PCA is a widely used technique that transforms the original features into a new set of uncorrelated features called principal components. These components are ordered in terms of the amount of variance they explain in the data. PCA finds linear combinations of the original features that maximize the variance, allowing for dimensionality reduction while retaining the most important information.

b. Independent Component Analysis (ICA): ICA aims to find a linear transformation of the original features in which the resulting components are statistically independent. Unlike PCA, which focuses on uncorrelated components, ICA seeks to extract underlying independent sources that may have generated the observed data. ICA is often used in signal processing and image analysis.

c. Manifold Learning: Manifold learning methods aim to discover the underlying structure or geometry of the data. These techniques create a low-dimensional representation of the data by preserving the relationships between the original features. Examples of manifold learning algorithms include t-SNE (t-Distributed Stochastic Neighbor Embedding) and Isomap.

d. Autoencoders: Autoencoders are neural network models that are trained to reconstruct the input data from a compressed latent space representation. The encoder part of the network learns a lower-dimensional representation of the original features, while the decoder part reconstructs the original data from the compressed representation. By training the autoencoder to minimize the reconstruction error, the model learns a compressed representation that captures the essential features of the data.

## 10. Explain Spatial Data Mining, Types. Also explain CART Algorithm.

Spatial Data Mining refers to the process of discovering patterns, relationships, and knowledge from spatial data. Spatial data includes geographic information such as maps, satellite imagery, GPS data, and other location-based data. Spatial data mining techniques combine traditional data mining methods with spatial analysis and visualization to extract meaningful insights from spatial datasets.

Types of Spatial Data Mining:

1. Spatial Clustering: This type of spatial data mining focuses on identifying groups or clusters of spatial objects that are similar in terms of their geographic proximity or

attribute values. Spatial clustering techniques help in identifying spatial patterns and spatially coherent regions.

2. Spatial Classification: Spatial classification aims to assign spatial objects into predefined classes or categories based on their attributes and spatial relationships. It combines traditional classification algorithms with spatial analysis to classify objects in a spatial context.

3. Spatial Association Rule Mining: This type of spatial data mining focuses on discovering interesting relationships or associations among spatial objects. It helps in identifying patterns such as "hotspots" or areas with high concentrations of specific attributes.

4. Spatial Outlier Detection: Spatial outlier detection techniques identify spatial objects that deviate significantly from the expected behavior or distribution. Outliers may indicate unusual events, anomalies, or errors in the dataset.

5. Spatial Regression: Spatial regression techniques explore relationships between spatial objects and their attributes. It helps in modeling and predicting attribute values based on spatial factors.

Now let's explain the CART (Classification and Regression Trees) algorithm with pseudo code:
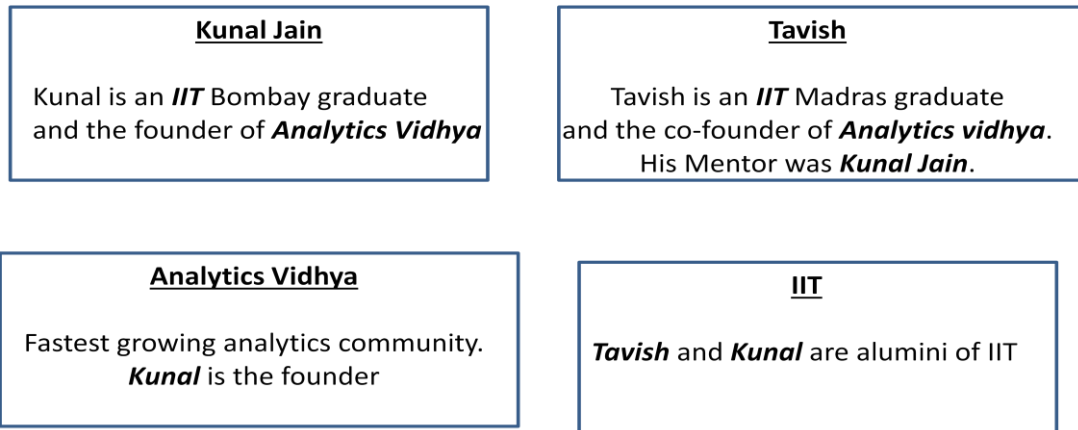
CART Algorithm:

Step 1: Start with the entire dataset D. Step 2: Select a target variable or attribute to predict. Step 3: If all instances in D belong to the same class or have the same attribute values, create a leaf node and label it with the corresponding class or attribute value. Return the leaf node. Step 4: For each attribute A, calculate the best split point that maximizes the separation of instances based on the target variable. Choose the attribute with the best split. Step 5: Create a decision node that tests the selected attribute. Step 6: Partition the dataset D into subsets D1, D2, ..., Dn based on the outcomes of the selected attribute test. Step 7: For each partition Di, repeat Steps 2-6 recursively until reaching a stopping criterion (e.g., a maximum tree depth or a minimum number of instances in a leaf node). Step 8: If the stopping criterion is met, create a leaf node for each partition Di and label them accordingly. Step 9: Return the decision tree.

The CART algorithm recursively partitions the dataset based on the selected attributes until it creates a decision tree that predicts the target variable with high accuracy. At each step, it chooses the attribute that best splits the data based on certain criteria, such as Gini index or information gain.

Note that the pseudo code provided is a simplified representation of the CART algorithm, and the actual implementation may involve additional details and considerations for handling specific data types and scenarios.

**Faculty:** Ms Jamuna S Murthy, Assistant Professor, Dept. of CSE, MSRIT, 2022-23

**Possible Questions:**

1. Imagine a web which has only 4 web pages, which are linked to each other. Each of the box below represents a web page. The words written in black and italics are the links between pages. Draw the directed graph of the web links and using page rank algorithm conclude the most important WebPage

| **Kunal Jain** | **Tavish** |
|---|---|
| Kunal is an *IIT* Bombay graduate and the founder of *Analytics Vidhya* | Tavish is an *IIT* Madras graduate and the co-founder of *Analytics vidhya*. His Mentor was *Kunal Jain*. |

| **Analytics Vidhya** | **IIT** |
|---|---|
| Fastest growing analytics community. *Kunal* is the founder | *Tavish* and *Kunal* are alumini of IIT |

2. What is dimensionality reduction in unsupervised learning and how does it differ from supervised learning?

3. How is subset selection used in dimensionality reduction techniques, and why is it important?

4. Can you explain the concept of Principal Component Analysis (PCA) and how it is utilized in dimensionality reduction?

5. Could you provide some numerical examples to illustrate the application of PCA in reducing the dimensionality of a dataset?

6. What is data mining, and how does it relate to different types of data?

7. What is the PageRank algorithm and how does it contribute to mining the World Wide Web?

8. Can you explain the process of text mining and its significance in extracting useful information from textual data?

9. How is mining time series data different from mining other types of data, and what are some techniques commonly used in this context?

10. What are ensemble methods, and how do they help in increasing the accuracy of data mining models?

11. Can you provide examples of ensemble methods used in practice to improve the performance of data mining algorithms?

**Faculty:** Ms Jamuna S Murthy, Assistant Professor, Dept. of CSE, MSRIT, 2022-23

12. How do unsupervised learning techniques, such as dimensionality reduction and data mining, contribute to knowledge discovery and decision-making processes in various industries?

13. What are the potential challenges or limitations associated with dimensionality reduction and data mining, and how can they be addressed?

14. How does the choice of data representation and feature selection impact the outcomes of dimensionality reduction and data mining tasks?

15. Can you discuss the ethical considerations and potential biases that may arise when applying dimensionality reduction and data mining techniques on large datasets?

16. In what ways can dimensionality reduction and data mining techniques be applied in real-world scenarios, such as customer segmentation, anomaly detection, or recommendation systems?

**Faculty:** Ms Jamuna S Murthy, Assistant Professor, Dept. of CSE, MSRIT, 2022-23