**Data Mining and Machine Learning**

**Lecture Notes – Module 3**

**Machine Learning Introduction:** Learning, Types of Machine Learning, Types of Machine Learning, Supervised Learning, The Machine Learning Process.

**Cluster Analysis:** Basic concepts and methods: Cluster Analysis, Partitioning methods, Hierarchical Methods, Evaluation of clustering.

**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. What Is Machine Learning? List out some Examples of Machine Learning Applications.

Machine Learning (ML) is a subset of artificial intelligence that involves the development of algorithms and statistical models that enable computer systems to automatically learn and improve from experience without being explicitly programmed. In other words, it is a method of teaching computers to learn from data and make predictions or take actions based on that learning.

Machine Learning relies on the analysis of large amounts of data to identify patterns, correlations, and insights that can be used to make informed decisions or predictions. It encompasses various techniques, such as supervised learning, unsupervised learning, reinforcement learning, and deep learning, each with its own set of algorithms and approaches.

Examples of Machine Learning applications can be found in numerous fields and industries. Here are a few notable ones:

1. Image Recognition: ML algorithms can be trained to recognize and classify images, enabling applications such as facial recognition, object detection, and self-driving cars.

2. Natural Language Processing (NLP): ML models can be used to analyze and understand human language, enabling applications like language translation, sentiment analysis, and chatbots.

3. Recommendation Systems: ML algorithms can analyze user preferences and behavior to provide personalized recommendations in areas such as movie streaming platforms, e-commerce websites, and music streaming services.

4. Fraud Detection: ML techniques can detect patterns and anomalies in large datasets, helping to identify fraudulent transactions, credit card fraud, or cybersecurity threats.

5. Medical Diagnosis: ML algorithms can assist in diagnosing diseases, analyzing medical images like X-rays and MRIs, and predicting patient outcomes based on historical data.

6. Financial Forecasting: ML models can analyze historical financial data to predict stock market trends, make investment decisions, and perform risk assessments.

7. Autonomous Vehicles: ML plays a crucial role in self-driving cars by enabling them to perceive and interpret their surroundings, make decisions, and navigate safely.

8. Virtual Assistants: ML powers virtual assistants like Siri, Alexa, and Google Assistant, allowing them to understand and respond to user queries, perform tasks,

and improve over time.

9. Customer Segmentation: ML techniques can analyze customer data to identify different segments and target specific marketing strategies based on individual preferences and behavior.

10. Energy Optimization: ML algorithms can optimize energy consumption by analyzing patterns and predicting demand, leading to more efficient energy management and cost savings.

**2. Explain the different types of Machine Learning.**

Here are the main types of machine learning:

1. Supervised Learning: Supervised learning is a type of machine learning where the algorithm learns from labeled training data. It involves mapping input data to desired output labels based on examples provided. The algorithm tries to generalize patterns and relationships in the data to make accurate predictions or classifications when presented with new, unseen data. Examples include image classification, spam filtering, and sentiment analysis.

2. Unsupervised Learning: Unsupervised learning deals with unlabeled data, where the algorithm aims to discover inherent patterns or structures in the data without explicit guidance. It explores the data's inherent structure, clustering similar data points together or reducing the dimensionality of the data. Unsupervised learning algorithms are often used for tasks like customer segmentation, anomaly detection, and recommender systems.

3. Semi-Supervised Learning: Semi-supervised learning is a combination of supervised and unsupervised learning. It is used when only a subset of the data is labeled, and the algorithm learns from both labeled and unlabeled data to make predictions or classifications. The labeled data helps guide the learning process, while the unlabeled data aids in capturing the underlying patterns. This approach is useful when labeling large amounts of data is expensive or time-consuming.

4. Reinforcement Learning: Reinforcement learning involves an agent that learns through interaction with an environment. The agent takes actions in the environment and receives feedback in the form of rewards or penalties based on its actions. The goal is to maximize the cumulative reward over time by learning which actions lead to the most favorable outcomes. Reinforcement learning is commonly used in robotics, gaming, and autonomous systems.

**3.  Explain Supervised Learning with Examples.**

Supervised learning is a machine learning approach where an algorithm learns a mapping between input data and corresponding output labels, based on a training dataset that is already labeled. The goal of supervised learning is to predict or classify new, unseen data accurately. It involves two main components: the input data (features) and the corresponding correct output labels.

To illustrate the concept of supervised learning, let's consider an example of email spam classification. Suppose we have a dataset of emails, each labeled as either spam or not spam. The input features can include various attributes of the email, such as the subject line, sender, and content. The corresponding output labels would be either "spam" or "not spam." The task of the supervised learning algorithm is to learn the patterns and characteristics in the input data that are indicative of spam emails, enabling it to correctly classify new, unseen emails as spam or not spam.

Now, let's explore two popular algorithms used in supervised learning:

1.    Decision Trees: Decision trees are a versatile and intuitive algorithm used for classification and regression tasks. They partition the input data based on different features and create a tree-like structure to make predictions. Each internal node represents a decision based on a particular feature, and each leaf node represents a predicted output label or value. For instance, in our email spam classification example, a decision tree might split the data based on the presence of certain keywords or the length of the subject line. Decision trees are easy to interpret and understand, but they can be prone to overfitting if not properly controlled.

2.    Support Vector Machines (SVM): SVM is a powerful algorithm commonly used for classification tasks. It aims to find an optimal hyperplane that separates the input data into different classes. The hyperplane is chosen such that the margin between the classes is maximized. SVMs can handle both linear and non-linear classification problems through the use of kernel functions. In our email spam classification example, SVM could learn a decision boundary that effectively separates spam and non-spam emails based on the chosen features. SVMs work well with small to medium-sized datasets and have been widely applied in various domains.

### 4. Explain Machine Learning Process in Detail

Here are the steps involved in the machine learning process:

1. **Data Collection:** The first step in the machine learning process is to gather relevant data for the problem at hand. This data can be obtained from various sources such as databases, APIs, or online repositories. The quality and quantity of the data collected play a crucial role in the success of the machine learning model.

2. **Data Preprocessing:** Once the data is collected, it needs to be preprocessed to ensure its quality and suitability for the model. This step involves tasks like data cleaning, handling missing values, removing outliers, and transforming the data into a format suitable for analysis.

3. **Feature Engineering:** Feature engineering involves selecting and creating relevant features from the available data that can contribute to the model's predictive power. This step requires domain knowledge and creativity to extract meaningful features that capture the underlying patterns in the data.

4. **Data Splitting:** In order to evaluate the performance of the machine learning model, the collected data is typically divided into training and testing sets. The training set is used to train the model, while the testing set is used to assess its performance on unseen data. It is important to maintain a proper balance between the two sets to avoid overfitting or underfitting.

5. **Model Selection:** The next step is to choose an appropriate machine learning algorithm or model that is best suited for the problem at hand. This decision depends on various factors such as the nature of the data, the type of problem (classification, regression, etc.), and the available computational resources.

6. **Model Training:** Once the model is selected, the training process begins. In this step, the model learns from the training data by adjusting its internal parameters to minimize the difference between the predicted outputs and the actual outputs. The optimization is usually achieved through techniques like gradient descent or stochastic gradient descent.

7. **Model Evaluation:** After the model is trained, it needs to be evaluated to assess its performance. This evaluation is done using the testing set that was set aside earlier. Common evaluation metrics include accuracy, precision, recall, and F1 score, depending on the problem type. The performance of the model helps determine if further adjustments or fine-tuning are required.

8. **Model Optimization:** If the model's performance is not satisfactory, optimization techniques can be applied to improve its accuracy or generalization. This can involve adjusting hyperparameters, changing the model architecture, or using ensemble methods to combine multiple models for better results.

9. **Model Deployment:** Once the model is optimized and meets the desired performance criteria, it can be deployed for real-world applications. This involves integrating the model into an existing system or creating a new application that utilizes its predictive capabilities. Deployment considerations include scalability, reliability, and security.

10. **Model Monitoring and Maintenance:** After deployment, it is crucial to continuously monitor the model's performance and re-evaluate it periodically. This ensures that the model continues to provide accurate predictions as the data distribution and patterns change over time. Maintenance tasks may include retraining the model with updated data or incorporating feedback from users to improve its performance.

**5. What is Cluster Analysis? Explain its different methods.**

Cluster analysis is a technique used in data mining and machine learning to identify groups or clusters within a dataset. It aims to partition data points into subsets that are similar within each cluster but dissimilar across clusters. Here are the basic concepts and methods of cluster analysis:

1. **Data representation:** Cluster analysis operates on a dataset that consists of a collection of objects or data points. These data points can be represented as vectors or multidimensional data.

2. **Similarity or distance measures:** The choice of a similarity or distance measure is crucial in cluster analysis. It quantifies the similarity or dissimilarity between pairs of data points. Common distance measures include Euclidean distance, Manhattan distance, and cosine similarity.

3. **Partitioning algorithms:** Partitioning algorithms aim to divide the dataset into non-overlapping clusters. One of the most well-known algorithms is the k-means algorithm, where k represents the desired number of clusters. It iteratively assigns data points to the nearest cluster centroid and updates the centroids until convergence.

4. **Hierarchical clustering:** Hierarchical clustering creates a tree-like structure of clusters, known as a dendrogram. It can be agglomerative (bottom-up) or divisive (top-down). Agglomerative clustering starts with each data point as a separate cluster and merges the most similar clusters iteratively until a single cluster is formed.

Divisive clustering starts with the entire dataset as one cluster and split it into smaller clusters recursively.

5. **Density-based clustering:** Density-based clustering algorithms identify dense regions of data points separated by sparser regions. The most popular density-based clustering algorithm is DBSCAN (Density-Based Spatial Clustering of Applications with Noise), which groups together data points that are close to each other and have a sufficient number of nearby neighbors.

6. **Model-based clustering:** Model-based clustering assumes that the data points are generated from a mixture of probability distributions. It attempts to find the best-fitting model that explains the data distribution and assigns data points to clusters based on the model parameters. The Gaussian Mixture Model (GMM) is a commonly used model-based clustering algorithm.

7. **Evaluation and validation:** Various metrics can be used to evaluate and validate the quality of clustering results. These metrics include the silhouette coefficient, Davies-Bouldin index, and Rand index. They assess the compactness and separation of clusters or compare the clustering results with known ground truth labels.

8. **Interpretation and visualization:** Once the clustering process is complete, interpreting and visualizing the results become important. Techniques such as scatter plots, heatmaps, and parallel coordinate plots can be used to understand the characteristics and differences among the clusters.

**6. What are partition clustering algorithms?**

Partitioning method: Partitioning a database $D$ of $n$ objects into a set of $k$ clusters, such that the sum of squared distances is minimized (where $c_i$ is the centroid or medoid of cluster $C_i$)

Given $k$, find a partition of $k$ *clusters* that optimizes the chosen partitioning criterion

Global optimal: exhaustively enumerate all partitions

Heuristic methods: *k-means* and *k-medoids* algorithms

- *k-means* (MacQueen'67, Lloyd'57/'82): Each cluster is represented by the center of the cluster

- *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

**7. Explain K-Means Clustering Algorithm with and example.**

K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if K=2, there will be two clusters, and for K=3, there will be three clusters, and so on.

It is an iterative algorithm that divides the unlabeled dataset into k different clusters in such a way that each dataset belongs only one group that has similar properties.

Backward Skip 10sPlay VideoForward Skip 10s

It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.
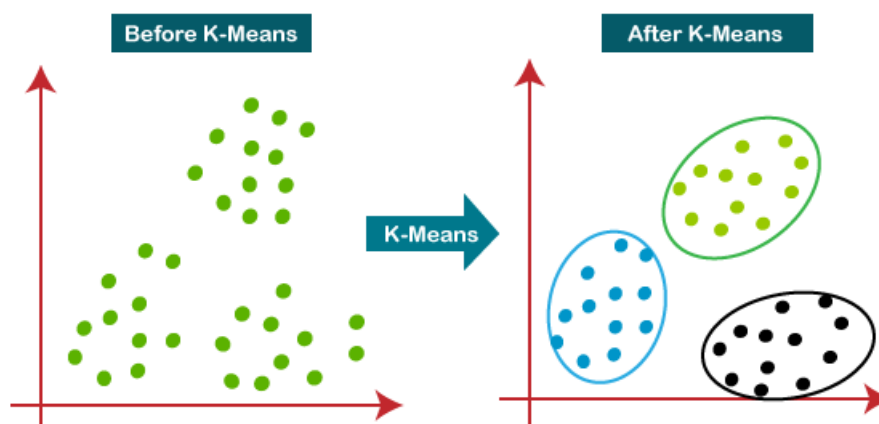
The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it does not find the best clusters. The value of k should be predetermined in this algorithm.

The k-means clustering algorithm mainly performs two tasks:

o   Determines the best value for K center points or centroids by an iterative process.
o   Assigns each data point to its closest k-center. Those data points which are near to the particular k-center, create a cluster.

Hence each cluster has datapoints with some commonalities, and it is away from other clusters.

The below diagram explains the working of the K-means Clustering Algorithm:

The working of the K-Means algorithm is explained in the below steps:

**Step-1:** Select the number K to decide the number of clusters.

**Step-2:** Select random K points or centroids. (It can be other from the input dataset).

**Step-3:** Assign each data point to their closest centroid, which will form the predefined K clusters.

**Step-4:** Calculate the variance and place a new centroid of each cluster.

**Step-5:** Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

**Step-6:** If any reassignment occurs, then go to step-4 else go to FINISH.

**Step-7**: The model is ready.

8. **Cluster the following eight points (with (x, y) representing locations) into three clusters: A1(2, 10), A2(2, 5), A3(8, 4), A4(5, 8), A5(7, 5), A6(6, 4), A7(1, 2), A8(4, 9). Initial cluster centers are: A1(2, 10), A4(5, 8) and A7(1, 2). Use K-Means Algorithm to find the three cluster centers after the second iteration.**

The distance function between two points a = (x1, y1) and b = (x2, y2) is defined as- P(a, b) = |x2 − x1| + |y2 − y1|.

**Iteration-01:**

- We calculate the distance of each point from each of the center of the three clusters.
- The distance is calculated by using the given distance function.

The following illustration shows the calculation of distance between point A1(2, 10) and each of the center of the three clusters-

**Calculating Distance Between A1(2, 10) and C1(2, 10)-**

P(A1, C1)

= |x2 − x1| + |y2 − y1|

= |2 − 2| + |10 − 10|

= 0

**Calculating Distance Between A1(2, 10) and C2(5, 8)-**

P(A1, C2)

= |x2 − x1| + |y2 − y1|

= |5 − 2| + |8 − 10|

= 3 + 2

= 5

**Calculating Distance Between A1(2, 10) and C3(1, 2)-**

P(A1, C3)

= |x2 − x1| + |y2 − y1|

= |1 − 2| + |2 − 10|

= 1 + 8

= 9

In the similar manner, we calculate the distance of other points from each of the center of the three clusters.

Next,
- We draw a table showing all the results.
- Using the table, we decide which point belongs to which cluster.
- The given point belongs to that cluster whose center is nearest to it.

| Given Points | Distance from center (2, 10) of Cluster-01 | Distance from center (5, 8) of Cluster-02 | Distance from center (1, 2) of Cluster-03 | Point belongs to Cluster |
|---|---|---|---|---|
| A1(2, 10) | 0 | 5 | 9 | C1 |
| A2(2, 5) | 5 | 6 | 4 | C3 |
| A3(8, 4) | 12 | 7 | 9 | C2 |
| A4(5, 8) | 5 | 0 | 10 | C2 |

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

| A5(7, 5) | 10 | 5 | 9 | C2 |
| A6(6, 4) | 10 | 5 | 7 | C2 |
| A7(1, 2) | 9 | 10 | 0 | C3 |
| A8(4, 9) | 3 | 2 | 10 | C2 |

From here, New clusters are-

**Cluster-01:**

First cluster contains points-
- A1(2, 10)

**Cluster-02:**

Second cluster contains points-
- A3(8, 4)
- A4(5, 8)
- A5(7, 5)
- A6(6, 4)
- A8(4, 9)

**Cluster-03:**

Third cluster contains points-
- A2(2, 5)
- A7(1, 2)

Now,
- We re-compute the new cluster clusters.
- The new cluster center is computed by taking mean of all the points contained in that cluster.

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

**For Cluster-01:**

- We have only one point A1(2, 10) in Cluster-01.
- So, cluster center remains the same.

**For Cluster-02:**

Center of Cluster-02
= ((8 + 5 + 7 + 6 + 4)/5, (4 + 8 + 5 + 4 + 9)/5)
= (6, 6)

**For Cluster-03:**

Center of Cluster-03
= ((2 + 1)/2, (5 + 2)/2)
= (1.5, 3.5)

This is completion of Iteration-01.

## Iteration-02:

- We calculate the distance of each point from each of the center of the three clusters.
- The distance is calculated by using the given distance function.

The following illustration shows the calculation of distance between point A1(2, 10) and each of the center of the three clusters-

**Calculating Distance Between A1(2, 10) and C1(2, 10)-**

P(A1, C1)
= |x2 − x1| + |y2 − y1|
= |2 − 2| + |10 − 10|
= 0

**Calculating Distance Between A1(2, 10) and C2(6, 6)-**

P(A1, C2)

= |x2 − x1| + |y2 − y1|

= |6 − 2| + |6 − 10|

= 4 + 4

= 8

**Calculating Distance Between A1(2, 10) and C3(1.5, 3.5)-**

P(A1, C3)

= |x2 − x1| + |y2 − y1|

= |1.5 − 2| + |3.5 − 10|

= 0.5 + 6.5

= 7

In the similar manner, we calculate the distance of other points from each of the center of the three clusters.

Next,

- We draw a table showing all the results.
- Using the table, we decide which point belongs to which cluster.
- The given point belongs to that cluster whose center is nearest to it.

| Given Points | Distance from center (2, 10) of Cluster-01 | Distance from center (6, 6) of Cluster-02 | Distance from center (1.5, 3.5) of Cluster-03 | Point belongs to Cluster |
|---|---|---|---|---|
| A1(2, 10) | 0 | 8 | 7 | C1 |
| A2(2, 5) | 5 | 5 | 2 | C3 |

| A3(8, 4) | 12 | 4 | 7 | C2 |
|----------|----|----|----|----|
| A4(5, 8) | 5 | 3 | 8 | C2 |
| A5(7, 5) | 10 | 2 | 7 | C2 |
| A6(6, 4) | 10 | 2 | 5 | C2 |
| A7(1, 2) | 9 | 9 | 2 | C3 |
| A8(4, 9) | 3 | 5 | 8 | C1 |

From here, New clusters are-

**Cluster-01:**

First cluster contains points-

- A1(2, 10)
- A8(4, 9)

**Cluster-02:**

Second cluster contains points-

- A3(8, 4)
- A4(5, 8)
- A5(7, 5)
- A6(6, 4)

**Cluster-03:**

Third cluster contains points-

- A2(2, 5)
- A7(1, 2)

Now,

- We re-compute the new cluster clusters.
- The new cluster center is computed by taking mean of all the points contained in that cluster.

**For Cluster-01:**

Center of Cluster-01
= ((2 + 4)/2, (10 + 9)/2)
= (3, 9.5)

**For Cluster-02:**

Center of Cluster-02
= ((8 + 5 + 7 + 6)/4, (4 + 8 + 5 + 4)/4)
= (6.5, 5.25)

**For Cluster-03:**

Center of Cluster-03
= ((2 + 1)/2, (5 + 2)/2)
= (1.5, 3.5)

This is completion of Iteration-02.

After second iteration, the center of the three clusters are-
- C1(3, 9.5)
- C2(6.5, 5.25)
- C3(1.5, 3.5)

9. **Explain K-Medoids with and Example.**

K-Medoids is a clustering algorithm that is used to partition a given dataset into K clusters. It is an extension of the popular K-Means algorithm, which is used to partition a dataset into K clusters based on the Euclidean distance between the data points and their centroids.

In K-Medoids, instead of calculating the centroid of each cluster, the algorithm chooses a

representative point, called a medoid, from each cluster. The medoid is the data point that has the smallest sum of distances to all other points in the same cluster. This makes the algorithm more robust to noise and outliers in the dataset.

Here's the pseudo-code for the K-Medoids algorithm:

1. Initialize: Choose K data points from the dataset to be the initial medoids

2. Repeat until convergence:

    a. Assign each data point to the closest medoid

    b. For each medoid m:

       i. Compute the total cost of swapping m with each non-medoid data point

       ii. Choose the non-medoid data point that results in the smallest cost, and make it the new medoid

3. Return the K clusters and their medoids

**Example: PAM Algorithm**

The Partitioning Around Medoids (PAM) algorithm is a clustering algorithm that was introduced by Kaufman and Rousseeuw in 1987. It is a variant of the K-medoids algorithm, which is a partitional clustering algorithm. PAM is commonly used for clustering analysis when the data set is small to moderate in size.

Here's a step-by-step explanation of the PAM algorithm:

1. Initialize: Select k points from the dataset as the initial medoids (representative points) randomly or using a heuristic method.

2. Assign: Assign each data point to the closest medoid based on a distance metric, typically using Euclidean distance or other dissimilarity measures.

3. Swap: For each medoid, consider swapping it with a non-medoid point and compute the total cost (sum of distances) of the resulting configuration.

4. Select the best swap: Choose the swap that minimizes the total cost. If the cost is reduced, update the medoid with the non-medoid point.

5. Repeat steps 2-4 until no more improvements can be made or a predefined number of iterations is reached.

6. Output: The final medoids represent the clusters, and each data point is assigned to the cluster represented by its closest medoid.

**10. Explain Hierarchical clustering methods AGNES (Agglomerative Nesting) and DIANA (Divisive Analysis).**

Hierarchical clustering is a class of clustering algorithms that organizes data points into a hierarchy of clusters. Two popular hierarchical clustering methods are AGNES (Agglomerative Nesting) and DIANA (Divisive Analysis).

1. **AGNES (Agglomerative Nesting):** AGNES is an agglomerative hierarchical clustering algorithm. It starts with each data point as an individual cluster and then iteratively merges the closest clusters based on a chosen distance metric until a termination condition is met.

Here are the steps involved in AGNES:

    i.    Initialization: Begin with each data point as a separate cluster.

    ii.    Distance calculation: Compute the proximity (distance) matrix between each pair of clusters. The proximity between two clusters can be measured using various metrics such as Euclidean distance, Manhattan distance, or others.

    iii.    Cluster merging: Find the two closest clusters based on the chosen proximity measure and merge them into a single cluster. Update the proximity matrix to reflect the new distances between the merged cluster and the remaining clusters.

    iv.    Repeat: Repeat steps 2 and 3 until a termination condition is met. This condition can be a specific number of desired clusters or a threshold distance value.

    v.    Hierarchical tree formation: The result is a hierarchy of clusters represented by a dendrogram, which shows the merging process and the distances at which clusters were merged.

2. **DIANA (Divisive Analysis):** DIANA is a divisive hierarchical clustering algorithm. Unlike AGNES, DIANA starts with a single cluster containing all data points and then recursively divides the clusters into smaller subclusters until a termination condition is satisfied.

Here are the steps involved in DIANA:

    i.    Initialization: Start with a single cluster containing all data points.

    ii.    Distance calculation: Compute the proximity (distance) matrix between each pair of data points.

    iii.    Cluster splitting: Find the data point or subset of data points that maximizes the dissimilarity within the cluster. Split the cluster into two subclusters based on this criterion.

    iv.    Repeat: Recursively apply the cluster splitting step to each newly formed

subcluster until a termination condition is met. This condition can be a specific number of desired clusters or a threshold dissimilarity value.

v. Hierarchical tree formation: The result is a dendrogram representing the hierarchy of clusters, showing the splitting process and the dissimilarities at which clusters were divided.

Both AGNES and DIANA have their strengths and weaknesses. AGNES tends to be computationally efficient and is suitable for larger datasets, but it may suffer from the "chaining effect" where the initial merging decisions cannot be undone. On the other hand, DIANA allows for more flexibility in terms of cluster shapes but can be computationally expensive, especially for large datasets.

## 11. Explain BIRCH algorithm in details.

BIRCH (Balanced Iterative Reducing and Clustering Using Hierarchies) is a popular hierarchical clustering algorithm that efficiently handles large datasets. It constructs a clustering hierarchy called the CF tree (Clustering Feature tree) and provides an approximate representation of the dataset. Here are the steps involved in the BIRCH algorithm:

1. Initialize the CF tree: Create an empty CF tree with the desired branching factor (B), maximum number of CF entries (M), and threshold value (T). The branching factor determines the maximum number of child nodes in each non-leaf node, and the maximum number of CF entries determines the capacity of each node. The threshold value is used to decide if a new entry should be inserted into an existing node or create a new node.

2. Reading data points: Read the input data points one by one from the dataset.

3. Inserting data points into the CF tree: For each data point, insert it into the CF tree by following these steps: a. Start at the root node of the CF tree. b. Compute the distance between the data point and each entry (sub-cluster) in the current node. c. If the distance to any sub-cluster is less than the threshold value (T), recursively traverse the corresponding child node. Otherwise, create a new sub-cluster in the current node. d. Update the statistics (centroid, diameter, number of points, etc.) of the affected sub-clusters during the traversal or creation process.

4. Condensing the CF tree: Periodically, check the CF tree for nodes with fewer than M entries. These nodes are considered "unbalanced." To condense the tree and reduce its size, follow these steps: a. Traverse the tree in a bottom-up manner. b. For each

unbalanced node encountered, merge its entries with its parent node. c. Update the statistics of the parent node accordingly.

5. Constructing the clustering hierarchy: After the CF tree has been constructed and condensed, it represents the approximate clustering of the dataset. You can traverse the CF tree to extract the desired number of clusters or construct a dendrogram to visualize the hierarchy.

6. Finalize the clustering: Based on the hierarchical structure obtained in the previous step, you can apply a suitable technique (e.g., cutting the dendrogram at a specific height) to obtain the final clustering result.

The BIRCH algorithm provides an efficient way to handle large datasets by constructing an approximate representation of the data in the CF tree. It reduces the overall computational complexity by condensing the tree and avoids repeated computations by storing the statistics of each sub-cluster.

## 12. Explain CHAMELEON (Hierarchical Clustering Using Dynamic Modeling) Algorithm.

CHAMELEON (Hierarchical Clustering Using Dynamic Modeling) is an algorithm used for hierarchical clustering, which groups similar data points together based on their distances. It aims to handle datasets with different densities and shapes by adapting the clustering process dynamically. Here are the steps involved in the CHAMELEON algorithm:

1. Input: The algorithm takes a dataset consisting of n data points as input. Each data point is represented by a set of d-dimensional feature vectors.

2. Similarity Matrix: A similarity matrix is constructed to represent the pairwise similarity between all data points in the dataset. Various distance measures can be used to calculate the similarity, such as Euclidean distance or cosine similarity.

3. Partitioning: The dataset is initially partitioned into a set of initial clusters. The number of initial clusters can be determined based on domain knowledge or by using an existing clustering algorithm like k-means.

4. Inter-Cluster Similarity: The inter-cluster similarity matrix is calculated to capture the similarity between clusters. It is based on the average similarity of all pairs of data points from different clusters.

5. Merge Step: The algorithm starts merging clusters based on their similarity. Initially, each cluster is considered as a separate component. The merging process is performed iteratively until a stopping criterion is met.

6. Dynamic Modeling: CHAMELEON utilizes a dynamic modeling technique to adaptively update the similarity matrix and inter-cluster similarity as clusters are merged. The algorithm takes into account both the distances between individual data points and the distances between clusters.

7. Cluster Similarity Matrix: As clusters are merged, a cluster similarity matrix is constructed to represent the similarity between clusters. This matrix is used to guide the merging process, favoring the merging of clusters that are more similar to each other.

8. Agglomerative Clustering: The algorithm uses an agglomerative clustering approach, where it repeatedly merges the most similar clusters based on their similarity scores until a termination condition is reached. The termination condition can be based on the desired number of clusters or a predefined threshold for similarity.

9. Output: The final output of the CHAMELEON algorithm is a hierarchical clustering structure, represented as a dendrogram. The dendrogram shows the hierarchical relationship between clusters, allowing different levels of granularity in the clustering results.

Overall, the CHAMELEON algorithm dynamically adapts the clustering process based on the local density and shape characteristics of the dataset. By considering both the distances between individual data points and the distances between clusters, it can handle datasets with different densities and shapes more effectively than traditional clustering algorithms.

**13. Explain Evaluation methods for clustering in detail.**

Evaluation methods for clustering aim to assess the quality and effectiveness of clustering algorithms. Here are ten points explaining different evaluation methods for clustering:

1. External Indexes: These methods compare the clustering results with known ground truth labels, such as class labels in classification tasks. Examples include the Rand Index and the Fowlkes-Mallows Index.

2. Internal Indexes: These methods evaluate the clustering based on the internal structure of the data, without using any external information. Silhouette Coefficient and Davies-Bouldin Index are commonly used internal evaluation measures.

3. Centroid-Based Evaluation: This method assesses the quality of clustering by measuring the distance between cluster centroids and the data points within the clusters. The lower the intra-cluster distance and the higher the inter-cluster distance, the better the clustering.

4. Connectivity-Based Evaluation: This method evaluates the clustering based on the

connectivity of data points within clusters. It focuses on finding compact and well-connected clusters. Examples include the Connectivity Index and the Dunn Index.

5. Distribution-Based Evaluation: These methods evaluate clustering based on the distribution of data points within clusters. The goal is to achieve clusters with similar shapes and sizes. One popular measure is the Hopkins statistic.

6. Stability-Based Evaluation: This approach assesses the stability of clustering results by considering the robustness of the algorithm to perturbations in the data or variations in the algorithm parameters. Stability Index and Variation of Information are commonly used stability measures.

7. Entropy-Based Evaluation: This method evaluates the clustering quality by measuring the information entropy of the clusters. It quantifies the homogeneity of data points within clusters. Lower entropy indicates more homogeneous clusters.

8. Compactness-Based Evaluation: This method focuses on measuring the compactness of clusters. It assesses how tightly grouped the data points are within each cluster. Compactness measures include Intra-Cluster Variance and Sum of Squared Errors (SSE).

9. Overlapping Evaluation: In some scenarios, clustering algorithms may produce overlapping clusters. Overlapping evaluation methods assess the quality of such clustering results. Jaccard Index and Fuzzy Overlapping Measure are examples of such measures.

10. Visual Evaluation: Sometimes, it is valuable to visually inspect the clustering results to evaluate their quality. Visualization techniques, such as scatter plots or heatmaps, can help identify patterns, separability, and coherence of clusters.

These evaluation methods provide different perspectives for assessing the quality of clustering results. It is often advisable to employ multiple evaluation measures to gain a comprehensive understanding of the clustering algorithm's performance.

**Possible Questions:**

**Machine Learning Introduction:**

1. What is machine learning and how does it differ from traditional programming?
2. Explain the concept of training data in machine learning.
3. What are the main goals of machine learning?
4. Discuss the relationship between machine learning and artificial intelligence.
5. Describe the three main components of a machine learning system.

**Types of Machine Learning:**

6. Compare and contrast supervised, unsupervised, and reinforcement learning.
7. Provide examples of real-world applications for each type of machine learning.
8. What is the fundamental difference between classification and regression in machine learning?
9. Discuss the concept of semi-supervised learning and its significance.
10. Explain the concept of transfer learning and its benefits in machine learning.

**Supervised Learning:**

11. What is supervised learning and how does it work?
12. Discuss the difference between the input features and target variable in supervised learning.
13. Provide examples of classification and regression problems in supervised learning.
14. What is the role of labeled data in supervised learning?
15. Explain the concept of overfitting in supervised learning and how to address it.

**The Machine Learning Process:**

16. Describe the steps involved in the machine learning process.
17. Discuss the importance of data preprocessing in the machine learning pipeline.
18. What are the common techniques for feature selection in machine learning?
19. Explain the concept of model evaluation and validation in machine learning.
20. Discuss the difference between underfitting and overfitting in the context of model performance.

**Cluster Analysis:**

21. What is cluster analysis and how does it differ from classification?
22. Discuss the main objectives and applications of cluster analysis.
23. Explain the concept of similarity measures in cluster analysis.
24. What are the challenges associated with cluster analysis?
25. Describe the concept of clustering tendency and its significance.

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

**Partitioning Clustering Methods and Types:**

26. Explain the basic idea behind partitioning clustering methods.

27. Discuss the strengths and weaknesses of the k-means algorithm.

28. What is the role of the number of clusters in partitioning clustering?

29. Compare and contrast k-means and k-medoids clustering algorithms.

30. Explain the concept of centroid initialization in partitioning clustering.

**Hierarchical Clustering Methods and Types:**

31. Describe the principles behind hierarchical clustering methods.

32. Discuss the difference between agglomerative and divisive hierarchical clustering.

33. What are the advantages and disadvantages of hierarchical clustering?

34. Explain the concept of dendrograms in hierarchical clustering.

35. Compare and contrast single-linkage and complete-linkage clustering methods.

**Evaluation of Clustering:**

36. What are the common evaluation metrics used for assessing clustering quality?

37. Discuss the concept of internal and external validation in clustering.

38. Explain the silhouette coefficient and its interpretation in clustering.

39. What are the limitations of evaluation metrics in clustering?

40. Describe the concept of stability-based evaluation in clustering.

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