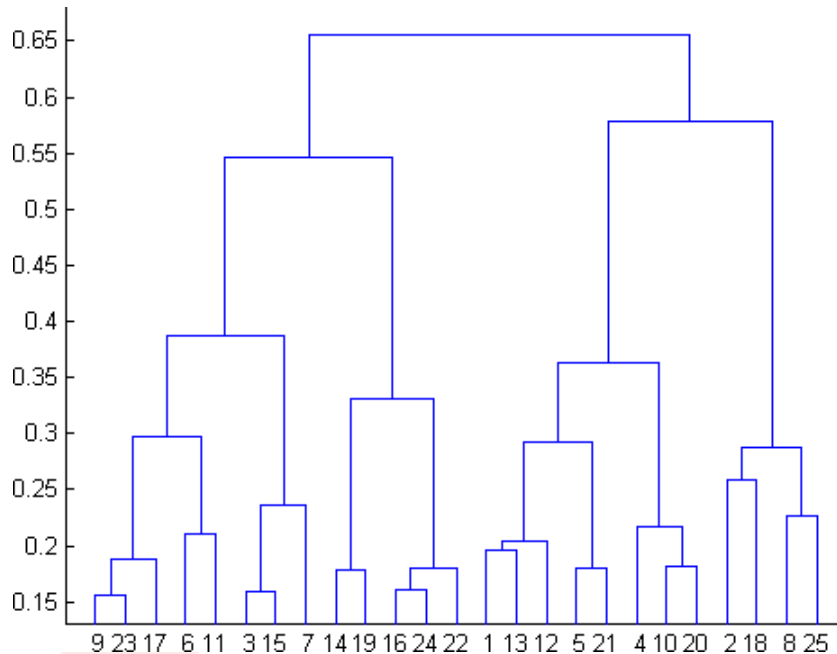**FLIP ROBO**

# MACHINE LEARNING

**Q1 to Q12 have only one correct answer. Choose the correct option to answer your question.**

1.  What is the most appropriate no. of clusters for the data points represented by the following dendrogram:



    a)  2
    b)  4
    c)  6
    d)  8
    Ans. 4

2.  In which of the following cases will K-Means clustering fail to give good results?
    1.  Data points with outliers
    2.  Data points with different densities
    3.  Data points with round shapes
    4.  Data points with non-convex shapes
    Options:
    a)  1 and 2
    b)  2 and 3
    c)  2 and 4
    d)  1, 2 and 4
    Ans.(D) 1, 2 and 4

3.  The most important part of_____is selecting the variables on which clustering is based.
    a)  interpreting and profiling clusters
    b)  selecting a clustering procedure
    c)  assessing the validity of clustering
    d)  formulating the clustering problem.
    Ans.(D) formulating the clustering problem.

4.  The most commonly used measure of similarity is the_____or its square.
    a)  Euclidean distance
    b)  city-block distance
    c)  Chebyshev's distance

# MACHINE LEARNING

    d) Manhattan distance
Ans. Euclidean distance

5. Which of the following is required by K-means clustering?
    a) Defined distance metric
    b) Number of clusters
    c) Initial guess as to cluster centroids
    d) All answers are correct
Ans. (d) All answers are correct

6. The goal of clustering is to-
    a) Divide the data points into groups
    b) Classify the data point into different classes
    c) Predict the output values of input data points
    d) All of the above
Ans.(a) Divide the data points into groups

7. Clustering is a-
    a) Supervised learning
    b) Unsupervised learning
    c) Reinforcement learning
    d) None
Ans.(b) Unsupervised learning

8. Which of the following clustering algorithms suffers from the problem of convergence at local optima?
    a) K- Means clustering
    b) Hierarchical clustering
    c) Diverse clustering
    d) All of the above
Ans(d) All of the above

9. Which version of the clustering algorithm is most sensitive to outliers?
    a) K-means clustering algorithm
    b) K-modes clustering algorithm
    c) K-medians clustering algorithm
    d) None
Ans.(a) K-means clustering algorithm

10. Which of the following is a bad characteristic of a dataset for clustering analysis-
    a) Data points with outliers
    b) Data points with different densities
    c) Data points with non-convex shapes
    d) All of the above
Ans(d)All of the above

11. For clustering, we do not require-
    a) Labeled data
    b) Unlabeled data
    c) Numerical data
    d) Categorical data
Ans(a) Labeled data

**Q13 to Q15 are subjective answers type questions, Answers them in their own words briefly.**

# MACHINE LEARNING

**12. How is cluster analysis calculated?**

**Ans.** The given data is divided into different groups by combining similar objects into a group. This group is nothing but a cluster. A cluster is nothing but a collection of similar data which is grouped together. Clustering Methods:

The clustering methods can be classified into the following categories:

- Partitioning Method
- Hierarchical Method
- Density-based Method
- Grid-Based Method
- Model-Based Method
- Constraint-based Method

**Partitioning Method:** It is used to make partitions on the data in order to form clusters. If "n" partitions are done on "p" objects of the database then each partition is represented by a cluster and n < p. The two conditions which need to be satisfied with this Partitioning Clustering Method are:

- One objective should only belong to only one group.
- There should be no group without even a single purpose.

In the partitioning method, there is one technique called iterative relocation, which means the object will be moved from one group to another to improve the partitioning

**Hierarchical Method:** In this method, a hierarchical decomposition of the given set of data objects is created. We can classify hierarchical methods and will be able to know the purpose of classification on the basis of how the hierarchical decomposition is formed. There are two types of approaches for the creation of hierarchical decomposition, they are:

- **Agglomerative Approach:** The agglomerative approach is also known as the bottom-up approach. Initially, the given data is divided into which objects form separate groups. Thereafter it keeps on merging the objects or the groups that are close to one another which means that they exhibit similar properties. This merging process continues until the termination condition holds.
- **Divisive Approach:** The divisive approach is also known as the top-down approach. In this approach, we would start with the data objects that are in the same cluster. The group of individual clusters is divided into small clusters by continuous iteration. The iteration continues until the condition of termination is met or until each cluster contains one object.

Once the group is split or merged then it can never be undone as it is a rigid method and is not so flexible. The two approaches which can be used to improve the Hierarchical Clustering Quality in Data Mining are: –

- One should carefully analyze the linkages of the object at every partitioning of hierarchical clustering.
- One can use a hierarchical agglomerative algorithm for the integration of hierarchical agglomeration. In this approach, first, the objects are grouped into micro-clusters. After grouping data objects into micro clusters, macro clustering is performed on the micro cluster.

**Density-Based Method:** The density-based method mainly focuses on density. In this method, the given cluster will keep on growing continuously as long as the density in the neighbourhood exceeds some threshold, i.e, for each data point within a given cluster. The radius of a given cluster has to contain at least a minimum number of points.

**Grid-Based Method:** In the Grid-Based method a grid is formed using the object together, i.e, the object space is quantized into a finite number of cells that form a grid structure. One of the major advantages of the grid-based method is fast processing time and it is dependent only on the number of cells in each dimension in the quantized space. The processing time for this method is much faster so it can save time.

**Model-Based Method:** In the model-based method, all the clusters are hypothesized in order to find the data which is best suited for the model. The clustering of the density function is used to locate the clusters for a given model. It reflects the spatial distribution of data points and also provides a way to automatically determine the number of clusters based on standard statistics, taking outlier or noise into account.Therefore it yields robust clustering methods.

**Constraint-Based Method:** The constraint-based clustering method is performed by the incorporation of application or user-oriented constraints. A constraint refers to the user expectation or the properties of the desired clustering results. Constraints provide us with an interactive way of communication with the clustering process. The user or the application requirement can specify constraints.

# MACHINE LEARNING

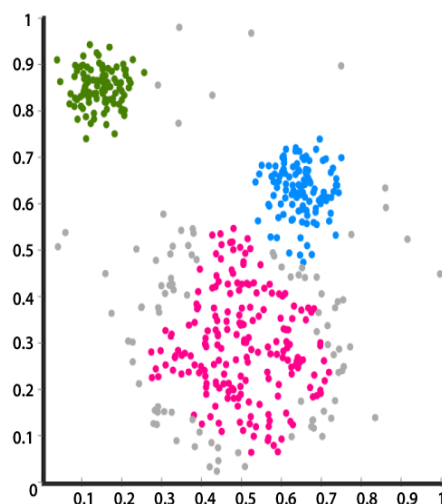13. How is cluster quality measured?

Ans. Measures for Quality of Clustering:

If all the data objects in the cluster are highly similar then the cluster has high quality. We can measure the quality of <u>Clustering</u> by using the Dissimilarity/Similarity metric in most situations. But there are some other methods to measure the Qualities of Good Clustering if the clusters are alike.

1. **Dissimilarity/Similarity metric:** The similarity between the clusters can be expressed in terms of a distance function, which is represented by $d(i, j)$. Distance functions are different for various data types and data variables. Distance function measure is different for continuous-valued variables, categorical variables, and vector variables. Distance function can be expressed as Euclidean distance, Mahalanobis distance, and Cosine distance for different types of data.

2. **Cluster completeness:** Cluster completeness is the essential parameter for good clustering, if any two data objects are having similar characteristics then they are assigned to the same category of the cluster according to ground truth. Cluster completeness is high if the objects are of the same category.

3. **Ragbag:** In some situations, there can be a few categories in which the objects of those categories cannot be merged with other objects. Then the quality of those cluster categories is measured by the Rag Bag method. According to the rag bag method, we should put the heterogeneous object into a rag bag category.

4. **Small cluster preservation:** If a small category of clustering is further split into small pieces, then those small pieces of cluster become noise to the entire clustering and thus it becomes difficult to identify that small category from the clustering. The small cluster preservation criterion states that are splitting a small category into pieces is not advisable and it further decreases the quality of clusters as the pieces of clusters are distinctive.

14. What is cluster analysis and its types?

Ans. Cluster analysis is a multivariate data mining technique whose goal is to groups objects (eg., products, respondents, or other entities) based on a set of user selected characteristics or attributes. It is the basic and most important step of data mining and a common technique for statistical data analysis.



Clusters should exhibit high internal homogeneity and high external heterogeneity.

# MACHINE LEARNING

## Types of Cluster Analysis

The clustering algorithm needs to be chosen experimentally unless there is a mathematical reason to choose one cluster method over another. It should be noted that an algorithm that works on a particular set of data will not work on another set of data. There are a number of different methods to perform cluster analysis. Some of them are,
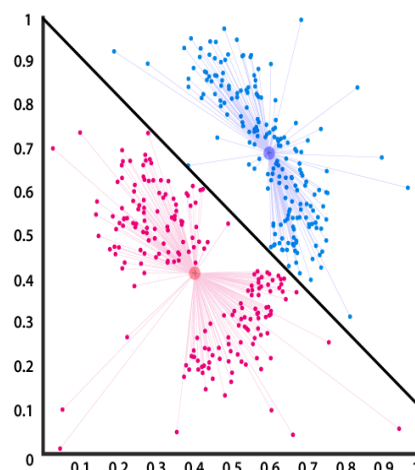
## Hierarchical Cluster Analysis

In this method, first, a cluster is made and then added to another cluster (the most similar and closest one) to form one single cluster. This process is repeated until all subjects are in one cluster. This particular method is known as **Agglomerative method**. Agglomerative clustering starts with single objects and starts grouping them into clusters.

**The divisive method** is another kind of Hierarchical method in which clustering starts with the complete data set and then starts dividing into partitions.
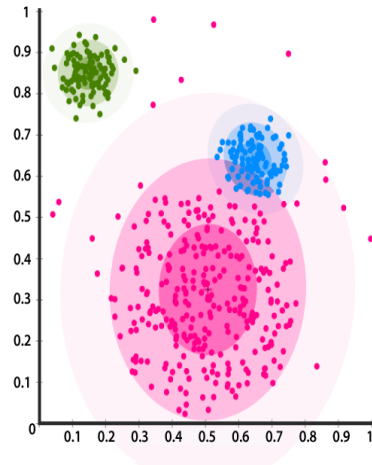
## Centroid-based Clustering

In this type of clustering, clusters are represented by a central entity, which may or may not be a part of the given data set. K-Means method of clustering is used in this method, where k are the cluster centers and objects are assigned to the nearest cluster centres.
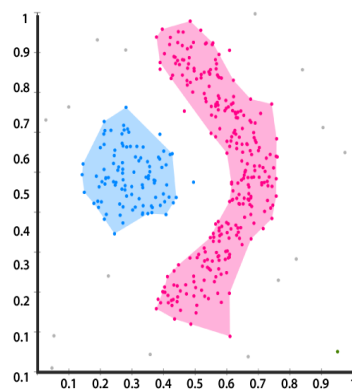


## Distribution-based Clustering

It is a type of clustering model closely related to statistics based on the modals of distribution. Objects that belong to the same distribution are put into a single cluster. This type of clustering can capture some complex properties of objects like correlation and dependence between attributes.

# MACHINE LEARNING



**Density-based Clustering**

In this type of clustering, clusters are defined by the areas of density that are higher than the remaining of the data set. Objects in sparse areas are usually required to separate clusters. The objects in these sparse points are usually noise and border points in the graph. The most popular method in this type of clustering is DBSCAN.

# MACHINE LEARNING