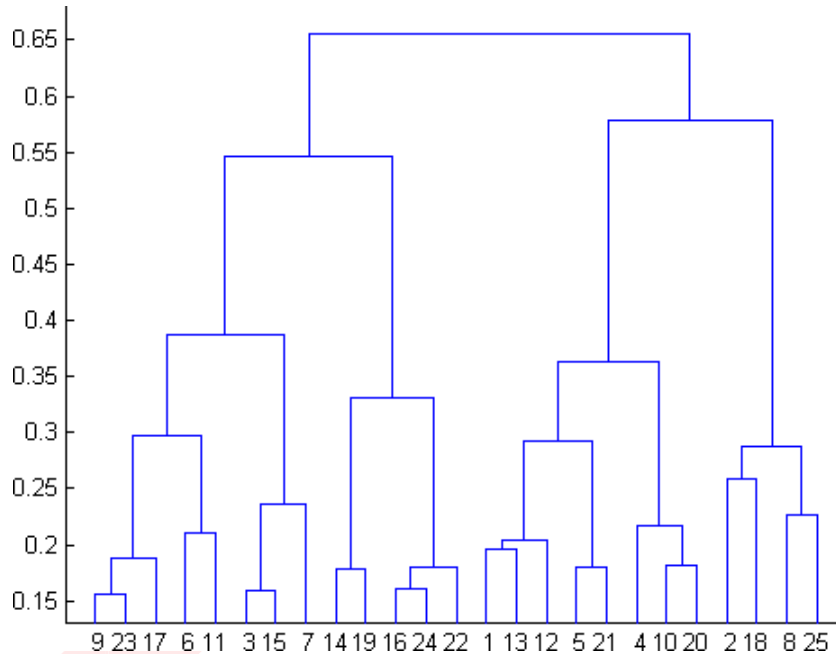**FLIP ROBO**

# 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

b)

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

    d)

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

    d)

4. The most commonly used measure of similarity is the _____ or its square.
    a) Euclidean distance

# MACHINE LEARNING

b) city-block distance
c) Chebyshev's distance
d) Manhattan distance

a)

# MACHINE LEARNING

5. ____ is a clustering procedure where all objects start out in one giant cluster. Clusters are formed by dividing this cluster into smaller and smaller clusters.
    a) Non-hierarchical clustering
    b) Divisive clustering
    c) Agglomerative clustering
    d) K-means clustering

    b)

6. 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

    d)

7. 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
       a)

8. Clustering is a-
    a) Supervised learning
    b) Unsupervised learning
    c) Reinforcement learning
    d) None
       b)

9. 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

       d)

10. 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

       a)

11. 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
       d)

12. For clustering, we do not require-
    a) Labeled data

# MACHINE LEARNING

b) Unlabeled data
c) Numerical data
d) Categorical data
a)

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

13. How is cluster analysis calculated?

This is calculated as the sum of squared distances between data points and the centers of the clusters they belong to. Inertia quantifies the within-cluster variation. Another popular metric is the silhouette coefficient, which attempts to summarize both within-cluster and between-cluster variation. The hierarchical cluster analysis follows three basic steps: 1) calculate the distances, 2) link the clusters, and 3) choose a solution by selecting the right number of clusters. First, we have to select the variables upon which we base our clusters.

14. How is cluster quality measured?

We have a few methods to choose from for measuring the quality of a clustering. In general, these methods can be categorized into two groups according to whether ground truth is available. Here, ground truth is the ideal clustering that is often built using human experts.
If ground truth is available, it can be used by extrinsic methods, which compare the clustering against the group truth and measure. If the ground truth is unavailable, we can use intrinsic methods, which evaluate the goodness of a clustering by considering how well the clusters are separated. Ground truth can be considered as supervision in the form of "cluster labels." Hence, extrinsic methods are also known as supervised methods, while intrinsic methods are unsupervised methods.

15. What is cluster analysis and its types?

Cluster analysis is a multivariate data mining technique whose goal is to groups objects ( 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, and it is used in many fields such as data compression, machine learning, pattern recognition, information retrieval etc.

Types of cluster analysis-

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.

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