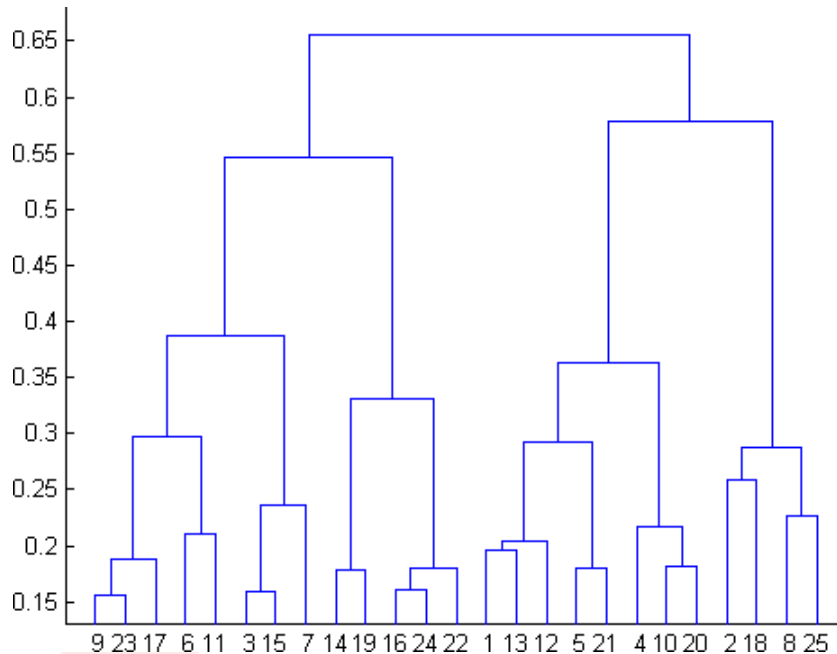


MACHINE LEARNING

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

Answer: (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

Answer : (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

Answer: (d)

MACHINE LEARNING

4. The most commonly used measure of similarity is the ____ or its square.
- Euclidean distance
 - city-block distance
 - Chebyshev's distance
 - Manhattan distance

Answer: (a)

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.
- Non-hierarchical clustering
 - Divisive clustering
 - Agglomerative clustering
 - K-means clustering

Answer: (b)

6. Which of the following is required by K-means clustering?
- Defined distance metric
 - Number of clusters
 - Initial guess as to cluster centroids
 - All answers are correct

Answer: (d)

7. The goal of clustering is to-
- Divide the data points into groups
 - Classify the data point into different classes
 - Predict the output values of input data points
 - All of the above

Answer: (a)

8. Clustering is a-
- Supervised learning
 - Unsupervised learning
 - Reinforcement learning
 - None

Answer: (b)

9. Which of the following clustering algorithms suffers from the problem of convergence at local optima?
- K- Means clustering
 - Hierarchical clustering
 - Diverse clustering
 - All of the above

Answer: (d)

10. Which version of the clustering algorithm is most sensitive to outliers?
- K-means clustering algorithm
 - K-modes clustering algorithm
 - K-medians clustering algorithm
 - None

Answer: (a)

11. Which of the following is a bad characteristic of a dataset for clustering analysis-

MACHINE LEARNING

- a) Data points with outliers
- b) Data points with different densities
- c) Data points with non-convex shapes
- d) All of the above

Answer: (d)

12. For clustering, we do not require-

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

Answer: (d)

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

13. How is cluster analysis calculated?

- 1) Copy the data into table
 - 2) select more than one variable
 - 3) select the no of cluster to calculate
- Clusters can be calculated various grouping methods
- graph-theoretical
- hierarchically
- partitioning
- optimizing

14. How is cluster quality measured?

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
-

MACHINE LEARNING

decreases the quality of clusters as the pieces of clusters are distinctive.

15. What is cluster analysis and its types?

Cluster Analysis is the process to find similar groups of objects in order to form clusters. It is an unsupervised machine learning-based algorithm that acts on unlabelled data. A group of data points would comprise together to form a cluster in which all the objects would belong to the same group. Cluster Analysis is the process to find similar groups of objects in order to form clusters. It is an unsupervised machine learning-based algorithm that acts on unlabelled data. A group of data points would comprise together to form a cluster in which all the objects would belong to the same group.

Hierarchical cluster Analysis

Centroid based clustering

Density based clustering
