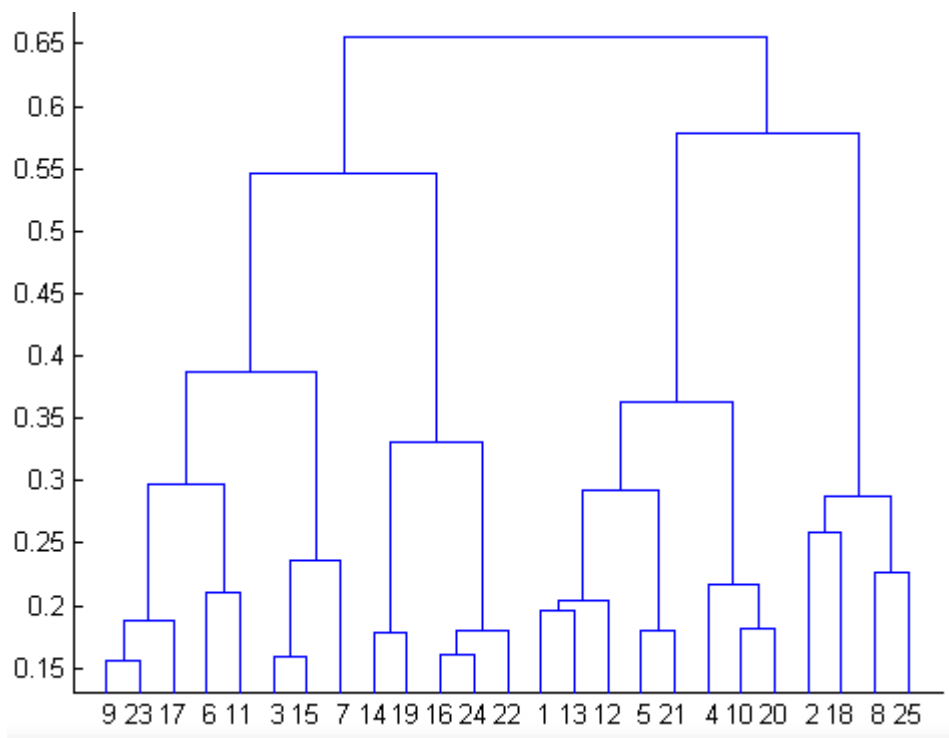


## MACHINE LEARNING

### ASSIGNMENT – 1

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

1. A

## MACHINE LEARNING

### ASSIGNMENT – 1

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

→ A

12. For clustering, we do not require-

- a) Labeled data
- 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?

→ Cluster analysis is an exploratory analysis that tries to identify structures within the data. Cluster analysis is also called segmentation analysis or taxonomy analysis. More specifically, it tries to identify homogenous groups of cases if the grouping is not previously known. Because it is exploratory, it does not make any distinction between dependent and independent variables. The different cluster analysis methods that SPSS offers can handle binary, nominal, ordinal, and scale (interval or ratio) data.

Cluster analysis is often used in conjunction with other analyses (such as discriminant analysis). The researcher must be able to interpret the cluster analysis based on their understanding of the data to determine if the results produced by the analysis are actually meaningful.

14. How is cluster quality measured?

→ We can measure the quality of Clustering 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.

Let us consider the clustering  $C1$ , which contains the sub-clusters  $s1$  and  $s2$ , where the members of the  $s1$  and  $s2$  cluster belong to the same category according to ground truth. Let us consider another clustering  $C2$  which is identical to  $C1$  but now  $s1$  and  $s2$  are merged into one cluster. Then, we define the clustering quality measure,  $Q$ , and according to cluster completeness  $C2$ , will have more cluster quality compared to the  $C1$  that is,  $Q(C2, C_g) > Q(C1, C_g)$ .

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.

Let us consider a clustering  $C1$  and a cluster  $C \in C1$  so that all objects in  $C$  belong to the same category of cluster  $C1$  except the object  $o$  according to ground truth. Consider a clustering  $C2$  which is identical to  $C1$  except that  $o$  is assigned to a cluster  $D$  which holds the objects of different categories. According to the ground truth, this situation is noisy and the quality of clustering is measured using the rag bag criteria. we define the clustering quality measure,  $Q$ , and according to rag bag method criteria  $C2$ , will have more cluster quality compared to the  $C1$  that is,  $Q(C2, C_g) > Q(C1, C_g)$ .

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. Suppose clustering  $C1$  has split into three clusters,  $C11 = \{d1, \dots, dn\}$ ,  $C12 = \{dn+1\}$ , and  $C13 = \{dn+2\}$ .

Let clustering  $C2$  also split into three clusters, namely  $C1 = \{d1, \dots, dn-1\}$ ,  $C2 = \{dn\}$ , and  $C3 = \{dn+1, dn+2\}$ . As  $C1$  splits the small category of objects and  $C2$  splits the big category which is

preferred according to the rule mentioned above the clustering quality measure  $Q$  should give a higher score to  $C_2$ , that is,  $Q(C_2, C_g) > Q(C_1, C_g)$ .

15. What is cluster analysis and its types?

→ Clustering or Cluster analysis is the method of grouping the entities based on similarities.

Defined as an unsupervised learning problem that aims to make training data with a given set of inputs but without any target values.

It is the process of finding similar structures in a set of unlabelled data to make it more understandable and manipulative.

Types of clustering are:

1. Connectivity-based Clustering (Hierarchical clustering)
2. Centroids-based Clustering (Partitioning methods)
3. Distribution-based Clustering
4. Density-based Clustering (Model-based methods)
5. Fuzzy Clustering
6. Constraint-based (Supervised Clustering)