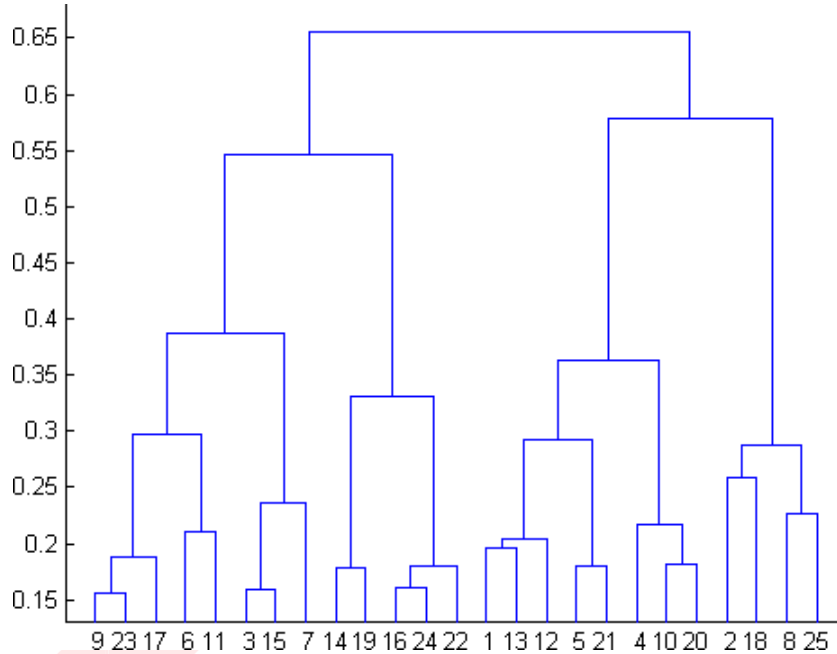


Q1 to Q12 have only one correct answer. Choose the correct option to answer your question.

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

Correct 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

Correct 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

Correct Answer: (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

Correct 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.
- a) Non-hierarchical clustering
 - b) Divisive clustering
 - c) **Agglomerative clustering**
 - d) K-means clustering

Correct Answer: (c)

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**

Correct Answer: (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

Correct Answer : (a)

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

Correct Answer: (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

FLIP ROBO

Correct Answer: (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

Correct Answer: (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

Correct Answer: (d)

12. For clustering, we do not require-

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

Correct Answer: (a)

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

13. How is cluster analysis calculated?

The cluster analysis is calculated following 3 basic steps:

- a) Calculate the distances
- b) Link the clusters
- c) Choose a solution by selecting the right number of clusters.

First we have to select the variables upon which we base our cluster. In the dialog window we add the math, reading and writing test to the list of variables. Since we want cluster cases we leave the rest of the tick marks on the default. In dialog box Statistics we can specify whether we want to output the proximity matrix and again we leave all settings on default. In the dialog box Plots we should add the Dendrogram. The Dendrogram will graphically show the clusters are merged and allow us to identify what the appropriate number of clusters. The dialog box Method allows us to specify the distance measures and the clustering method. First we need to define the correct distance measure. For interval

data, the most common is Square Euclidian Distance. It is based on Euclidian distance between two observation which is the squared root of the sum of the squared distance. Since the Euclidian Distance squared it increases the importance of large distances, while weakening the importance of small distances.

Next we have to choose the Cluster Method. Typically, choices are between groups linkage (distance between clusters is the average distance of all data points within these cluster), nearest neighbor (single linkage; distance between clusters is the smallest distance between two data points). Single linkage works best with long chains of clusters while complete linkage works best with dense blobs of clusters. Between groups linkage works with both cluster types. So single linkage is recommended. It helps in identifying outliers. After excluding these outliers we can move onto ward's method. The last step is Standardization. If the variable have different scales and means we want to standardize either to Z scores or by centering scales.

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.

Extrinsic Methods

When the ground truth is available, we can compare it with a clustering to assess the clustering. Thus, the core task in extrinsic methods is to assign a score, $Q(C, C_g)$, to a clustering, C , given the ground truth, C_g . Whether an extrinsic method is effective largely depends on the measure, Q , it uses.

Intrinsic Methods

When the ground truth of a data set is not available, we have to use an intrinsic method to assess the clustering quality. In general, intrinsic methods evaluate a clustering by examining how well the clusters are separated and how compact the clusters are. Many intrinsic methods have the advantage of a similarity metric between objects in the data set.

4. What is cluster analysis and its types?

Cluster analysis is the task of grouping a set of data points in such a way that they can be characterized by their relevancy to one another. These techniques create clusters that allow us to understand how our data is related. The most common applications of cluster analysis in a business setting is to segment customers or activities.

Cluster analysis is a class of techniques that are used to classify objects or cases into relative groups called clusters. Cluster analysis is also called classification analysis or numerical taxonomy. In cluster analysis, there is no prior information about the group or cluster membership for any of the objects. Cluster Analysis has been used in marketing for various purposes. Segmentation of consumers in cluster analysis is used on the basis of benefits sought from the purchase of the product. It can be used to identify homogeneous groups of buyers.

Cluster analysis involves formulating a problem, selecting a distance measure, selecting a clustering procedure, deciding the number of clusters, interpreting the profile clusters and finally, assessing the validity of clustering. The variables on which the cluster analysis is to be done should be selected by keeping past research in mind. It should also be selected by theory, the hypotheses being tested, and the judgment of the researcher. An appropriate measure of distance or similarity should be selected; the most commonly used measure is the Euclidean distance or its square.

Clustering procedures in cluster analysis may be hierarchical, non-hierarchical, or a two-step procedure. A hierarchical procedure in cluster analysis is characterized by the development of a tree like structure. A hierarchical procedure can be agglomerative or divisive. Agglomerative methods in cluster analysis consist of linkage methods, variance methods, and centroid methods. Linkage methods in cluster analysis are comprised of single linkage, complete linkage, and average linkage.

The non-hierarchical methods in cluster analysis are frequently referred to as K means clustering. The two-step procedure can automatically determine the optimal number of clusters by comparing the values of model choice criteria across different clustering solutions. The choice of clustering procedure and the choice of distance measure are interrelated. The relative sizes of clusters in cluster analysis should be meaningful. The clusters should be interpreted in terms of cluster centroids.

Types of Cluster Analysis

There are four basic types of cluster analysis used in data science. These types are

1. Centroid Clustering
 2. Density Clustering
 3. Distribution Clustering
 4. Connectivity Clustering.
-

Centroid Clustering

This is one of the more common methodologies used in cluster analysis. In centroid cluster analysis you choose the number of clusters that you want to classify. For example, if you're a pet store owner you may choose to segment your customer list by people who bought dog and/or cat products. The algorithm will start by randomly selecting centroids (cluster centers) to group the data points into the two pre-defined clusters. A line is then drawn separating the data points into the two clusters based on their proximity to the centroids. The algorithm will then reposition the centroid relative to all the points within each cluster. The centroids and points in a cluster will adjust through all iterations, resulting in optimized clusters. The result of this analysis is the segmentation of your data into the two clusters. In this example, the data set will be segmented into customers who own dogs and cats.

Density Clustering

Density clustering groups data points by how densely populated they are. To group closely related data points, this algorithm leverages the understanding that the more dense the data points...the more related they are. To determine this, the algorithm will select a random point then start measuring the distance between each point around it. For most density algorithms a predetermined distance between data points is selected to benchmark how closely points need to be to one another to be considered related.. Then, the algorithm will identify all other points that are within the allowed distance of relevance. This process will continue to iterate by selecting different random data points to start with until the best clusters can be identified.

Distribution Clustering

Distribution clustering identifies the probability that a point belongs to a cluster. Around each possible centroid .The algorithm defines the density distributions for each cluster, quantifying the probability of belonging based on those distributions The algorithm optimizes the characteristics of the distributions to best represent the data.

These maps look a lot like targets at an archery range. In the event that a data point hits the bulls eye on the map, then the probability of that person/object belonging to that cluster is 100%.

Each ring around the bulls eye represents lessening percentage or certainty.

Distribution clustering is a great technique to assign outliers to clusters, where as density clustering will not assign an outlier to a cluster.

Connectivity Clustering

Unlike the other three techniques of clustering analysis reviewed above, connectivity clustering initially recognizes each data point as its own cluster. The primary premise of this technique is

that points closer to each other are more related. The iterative process of this algorithm is to continually incorporate a data point or group of data points with other data points and/or groups until all points are engulfed into one big cluster. The critical input for this type of algorithm is determining where to stop the grouping from getting bigger.