

## ANSWERS:

- 1) b) 4
- 2) d) 1, 2 and 4
- 3) d) formulating the clustering problem
- 4) a) Euclidean distance
- 5) b) Divisive clustering
- 6) d) all answers are correct
- 7) a) divide the data points into groups
- 8) b) unsupervised learning
- 9) a) K-means clustering
- 10) a) K-means clustering algorithm
- 11) d) all of the above
- 12) a) Labeled data

13)

Three methods for the cluster analysis Calculation: K-Means Cluster, Hierarchical Cluster, and Two-Step Cluster.

**K-means cluster** is a method to quickly cluster large data sets. The researcher define the number of clusters in advance. This is useful to test different models with a different assumed number of clusters.

**Hierarchical cluster** is the most common method. It generates a series of models with cluster solutions from 1 (all cases in one cluster) to n (each case is an individual cluster). Hierarchical cluster also works with variables as opposed to cases; it can cluster variables together in a manner somewhat similar to factor analysis. In addition, hierarchical cluster analysis can handle nominal, ordinal, and scale data; however it is not recommended to mix different levels of measurement.

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.

**Two-step cluster** analysis identifies groupings by running pre-clustering first and then by running hierarchical methods. Because it uses a quick cluster algorithm upfront, it can handle large data sets that would take a long time to compute with hierarchical cluster methods. In this respect, it is a combination of the previous two approaches. Two-step clustering can handle scale and ordinal data in the same model, and it automatically selects the number of clusters.

## 14)

If all the data objects in the cluster are highly similar then the cluster has high quality. 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.

### **Extrinsic methods (If ground truth is available) :**

**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  $C_1$ , which contains the sub-clusters  $s_1$  and  $s_2$ , where the members of the  $s_1$  and  $s_2$  cluster belong to the same category according to ground truth. Let us consider another clustering  $C_2$  which is identical to  $C_1$  but now  $s_1$  and  $s_2$  are merged into one cluster. Then, we define the clustering quality measure,  $Q$ , and according to cluster completeness  $C_2$ , will have more cluster quality compared to the  $C_1$  that is,  $Q(C_2, C_g) > Q(C_1, C_g)$ .

**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  $C_1$  and a cluster  $C \in C_1$  so that all objects in  $C$  belong to the same category of cluster  $C_1$  except the object  $o$  according to ground truth. Consider a clustering  $C_2$  which is identical to  $C_1$  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  $C_2$ , will have more cluster quality compared to the  $C_1$  that is,  $Q(C_2, C_g) > Q(C_1, C_g)$ .

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 C2, that is,  $Q(C2, Cg) > Q(C1, Cg)$ .

Many clustering quality measures satisfy some of these four criteria. Here, we introduce the BCubed precision and recall metrics, which satisfy all four criteria.

**Intrinsic Methods(If ground truth is unavailable)** :To measure the quality of a clustering, we can also use the average silhouette coefficient value of all objects in the data set.

## 15)

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, and it is used in many fields such as data compression, machine learning, pattern recognition, information retrieval etc.

### 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.

### **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.