**Q1 to Q11 have only one correct answer. Choose the correct option to answer your question**.

1. Movie Recommendation systems are an example of:

i) Classification

ii) Clustering

iii) Regression

Options:

a) 2 Only

b) 1 and 2

c) 1 and 3

d) 2 and 3

**Ans: d**

2. Sentiment Analysis is an example of:

i) Regression

ii) Classification

iii) Clustering

iv) Reinforcement

Options:

a) 1 Only

b) 1 and 2

c) 1 and 3

d) 1, 2 and 4

**Ans: d**

3. Can decision trees be used for performing clustering?

a) True

b) False

**Ans: a**

4. Which of the following is the most appropriate strategy for data cleaning before performing clustering

analysis, given less than desirable number of data points:

i) Capping and flooring of variables

ii) Removal of outliers

Options:

a) 1 only

b) 2 only

c) 1 and 2

d) None of the above

**Ans: a**

5. What is the minimum no. of variables/ features required to perform clustering?

a) 0

b) 1

c) 2

d) 3

**Ans: b**

6. For two runs of K-Mean clustering is it expected to get same clustering results?

a) Yes

b) No

**Ans: b**

7. Is it possible that Assignment of observations to clusters does not change between successive iterations in K-Means?

a) Yes

b) No

c) Can't say

d) None of these

**Ans: a**

8. Which of the following can act as possible termination conditions in K-Means?

i) For a fixed number of iterations.

ii) Assignment of observations to clusters does not change between iterations. Except for cases with a bad local minimum.

iii) Centroids do not change between successive iterations.

iv) Terminate when RSS falls below a threshold.

Options:

a) 1, 3 and 4

b) 1, 2 and 3

c) 1, 2 and 4

d) All of the above

**Ans: d**

9. Which of the following algorithms is most sensitive to outliers?

a) K-means clustering algorithm

b) K-medians clustering algorithm

c) K-modes clustering algorithm

d) K-medoids clustering algorithm

**Ans: a**

10. How can Clustering (Unsupervised Learning) be used to improve the accuracy of Linear Regression

model (Supervised Learning):

i) Creating different models for different cluster groups.

ii) Creating an input feature for cluster ids as an ordinal variable.

iii) Creating an input feature for cluster centroids as a continuous variable.

iv) Creating an input feature for cluster size as a continuous variable.

Options:

a) 1 only

b) 2 only

c) 3 and 4

d) All of the above

**Ans: d**

11. What could be the possible reason(s) for producing two different dendrograms using agglomerative

clustering algorithms for the same dataset?

a) Proximity function used

b) of data points used

c) of variables used

d) All of the above

**Ans: d**

**Q12 to Q14 are subjective answers type questions, Answers them in their own words briefly**

**12. Is K sensitive to outliers?**

**Ans:** Yes, Outliers can have a significant impact on the result of k-mean clustering. This is because the algorithm relies on minimizing the within cluster sum of square and outliers can often increase this value. As result, it is often recommended to remove outliers from your dataset before running K-mean clustering.

The k-means algorithm updates the cluster centers by taking the average of all the data points that are closer to each cluster center. When all the points are packed nicely together, the average makes sense. However, when you have outliers, this can affect the average calculation of the whole cluster. As a result, this will push your cluster center closer to the outliers.

**13. Why is K means better?**

**Ans:** There are following reasons for K-mean is better:

1. Easy to understand and implement.
2. Computationally efficient for both training and prediction.
3. Guaranteed convergence.
4. It is scalable to a huge data set and also faster to large datasets.
5. it adapts the new examples very frequently.
6. Generalization of clusters for different shapes and sizes.

**14. Is K means a deterministic algorithm?**

**Ans: No,** one of the significant drawbacks of K-Means is its non-deterministic nature. K-Means starts with a random set of data points as initial centroids. This random selection influences the quality of the resulting clusters. Besides, each run of the algorithm for the same dataset may yield a different output.