

WORKSHEET 3 - MACHINE LEARNING

1. Which of the following is an application of clustering?

- a. Biological network analysis
- b. Market trend prediction
- c. Topic modeling
- d. All of the above

ANS.B

2. On which data type, we cannot perform cluster analysis?

- a. Time series data
- b. Text data
- c. Multimedia data
- d. None

ANS. D

3. Netflix's movie recommendation system uses?

- a. Supervised learning
- b. Unsupervised learning
- c. Reinforcement learning and Unsupervised learning
- d. All of the above

ANS.D

4. The final output of Hierarchical clustering is?

- a. The number of cluster centroids
- b. The tree representing how close the data points are to each other
- c. A map defining the similar data points into individual groups
- d. All of the above

ANS.B

5. Which of the step is not required for K-means clustering?

- a. A distance metric
- b. Initial number of clusters
- c. Initial guess as to cluster centroids
- d. None

ANS. D

6. Which of the following is wrong?

- a. k-means clustering is a vector quantization method
- b. k-means clustering tries to group n observations into k clusters
- c. k-nearest neighbour is same as k-means
- d. None

ANS. C

7. Which of the following metrics, do we have for finding dissimilarity between two clusters in hierarchical clustering?

- i. Single link
- ii. Complete link
- iii. Average-link Options:
 - a. 1 and 2
 - b. 1 and 3
 - c. 2 and 3
 - d. 1, 2 and 3

ANS.D

8. Which of the following are true?

- i. Clustering analysis is negatively affected by multicollinearity of features
- ii. Clustering analysis is negatively affected by heteroscedasticity Options:
 - a. 1 only
 - b. 2 only
 - c. 1 and 2
 - d. None of them

ANS. B

9. In the figure above, if you draw a horizontal line on y-axis for $y=2$. What will be the number of clusters formed?

- a. 2
- b. 4
- c. 3
- d. 5

ANS. A

10. For which of the following tasks might clustering be a suitable approach?

- a. Given sales data from a large number of products in a supermarket, estimate future sales for each of these products.
- b. Given a database of information about your users, automatically group them into different market segments.
- c. Predicting whether stock price of a company will increase tomorrow.
- d. Given historical weather records, predict if tomorrow's weather will be sunny or rainy.

ANS.A

11. Given, six points with the following attributes:

Which of the following clustering representations and dendrogram depicts the use of MIN or Single link proximity function in hierarchical clustering:

ANS.A

12. Given, six points with the following attributes:

Which of the following clustering representations and dendrogram depicts the use of MAX or Complete link proximity function in hierarchical clustering.

ANS.B

13. What is the importance of clustering?

ANS: Clustering is the method of identifying similar groups of data in a dataset. Clustering is important in data analysis and data mining applications. It is the task of grouping a set of objects so that objects in the same group are more like each other than to those in other groups. is the task of dividing the population or data points into several groups such that data points in the same groups are more similar to other data points in the same group and dissimilar to the data points in other groups. It is basically a collection of objects based on similarity and dissimilarity between them Clustering is very much important as it determines the intrinsic grouping among the unlabelled data present. There are no criteria for good clustering. It depends on the user, what is the criteria they may use which satisfy their need. For instance, we could be interested in finding representatives for homogeneous groups (data reduction), in finding "natural clusters" and describe their unknown

properties (“natural” data types), in finding useful and suitable groupings (“useful” data classes) or in finding unusual data objects (outlier detection). This algorithm must make some assumptions that constitute the similarity of points and each assumption make different and equally valid clusters.

14. How can I improve my clustering performance?

ANS: There are two important elements in improving the quality of clustering: improving the weights of the features in a document vector and creating a more appropriate distance measure. A good weighting technique can promote the good features of an object, and an appropriate distance measure can help bring similar features together. The next two sections explain how you can create custom feature-selection and distance-measurement classes.

1 Improving document vector generation: - A good document vector has the right kind of features, with higher weights assigned to the more important ones. In text data, there are two ways to improve the quality of a document vector: by removing noise and using a good weighting technique.

2 Custom distance measure: - If the vectors are of the highest quality, the biggest improvement in cluster quality comes from the choice of an appropriate distance measure. We’ve seen that the cosine distance is a good distance measure for clustering text documents. To illustrate the power of a custom distance measure, we create a different form of the cosine distance measure that exaggerates distances: it makes big distances bigger and small distances smaller.