

Amin Fadaee

Machine Learning HW5

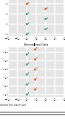
Python 3.5, Scikit-Learning, Matplotlib, Pandas, Numpy, Seaborn, sklearn Documentation, RStudio, JupyterLab, Jupyter

Problem 1

- **Distance** This is the most intuitive way to represent a cluster and has the least distortion within points inside cluster.
- **Median** Using mode is a simple way to represent a cluster as it is not influenced by the structure and the space of the points directly and also can be used in non-euclidean space as well.
- **Mean** Due to its robustness using median is most useful when there are outliers present in the data.

Problem 2

It depends with and without normalization will have **different** results. This is due to the fact that after normalization, the previous distances wouldn't be preserved and there is no guarantee that we will have the same clustering results. Here is an example with and demonstrating this fact.



As can be seen the obtained clusters are different for each case.

Problem 3

As 3 of the criteria have the same complexity if we want to find the linkage between two clusters we would have to find the distance between each point inside the cluster (because points forming points in the second cluster is points which results in the same complexity). Using the obtained distances d_{ij} the formula for deriving the linkages are as following:

$$\begin{aligned} \text{Single} &= \min(d_{ij}) \\ \text{Complete} &= \max(d_{ij}) \\ \text{Average} &= \frac{\sum d_{ij}}{n \times n} \end{aligned}$$

Single and Complete linkage are prone to outliers. If a point deviates from the majority of the points in cluster it can be considered a point too close or too far from the other clusters, however average linkage is robust to outliers.

Problem 4

The hyperparameter of DBSCAN is ϵ (radius) and the min_clusters . If we want to derive too large clusters a high value is needed and in the other hand by choosing a smaller value the number of clusters would increase. From other perspective if we know in advance the density of the clusters we can choose the min_clusters properly. A low min_clusters would result in low density clusters and a lot of dense cluster would require higher min_clusters .

In absence of strong prior knowledge of the clusters structure we can utilize the split score and validation set error as a method for finding the hyperparameters.

Problem 5

- a. The underlying this problem can be solved by means of $n \times n$ size the implementation of sklearn and DBSCAN respectively in sklearn `agg` and `min_samples`.

Here are the results for each k :

- $k=1$: 0.80101
- $k=2$: 0.79836
- $k=3$: 0.80101
- $k=4$: 0.80000

As can be seen $k=3$ would result in the best DB index. Here is the plot of the data clustered by 3 clusters with $\epsilon=0.5$.



- b. In the data is 2 dimensional average is good name of what could be good hyperparameters for the above algorithms. A good result getting because the shape of data is 2 clusters with 2 clusters. For achieving such clustering, $\epsilon=0.5$ and $\text{min_clusters}=3$ seems like a good choice. Here is a plot showing different clustering parameters.



- c. Using exploratory analysis and getting help from the data we can deduce that there are 4 clusters along with 2 outliers which cluster captured perfectly.

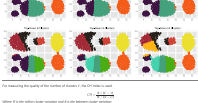
Problem 6

For the top-down clustering 3 clusters with $\epsilon=0.5$ is used. Also for finding the cluster to divide we pick a cluster with the highest average dissimilarity given by the formula below:

$$\text{dissim} = \sum_{i,j \in C_k} (x_i - x_j)^2$$

Where C_k is the cluster in cluster k .

Here are the results of the Top-Down Clustering:



For measuring the quality of the number of clusters k the CH index is used:

$$CH_k = \frac{W_k/B_k - W_{k+1}/B_{k+1}}{W_k/B_k}$$

Where W_k is the within-cluster variation and B_k is the between-cluster variation:

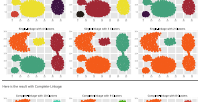
$$W_k = \sum_{i,j \in C_k} (x_i - x_j)^2$$

$$B_k = \sum_{i \in C_k} (x_i - \bar{x})^2$$

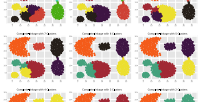
n_k is the number of point belonging to cluster k and \bar{x} is the mean of all data.

A complete description of the quality of the algorithms based on CH is discussed at the end of this problem.

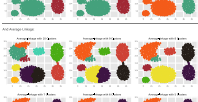
Using Single-linkage bottom-up clustering, we derived the following results:



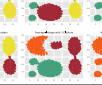
Here is the result with Complete Linkage:



And Average Linkage:



Using CH index we can compare the results of the above clustering algorithms following graph:

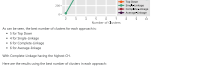


As can be seen, the best number of clusters for each approach:

- 3 for Top-Down
- 4 for Single-linkage
- 5 for Complete-linkage
- 6 for Average-linkage

With Complete Linkage having the highest CH.

Here are the results using the best number of clusters inside approach:



Because for this problem we are interested in top and bottom-up approach.