



University of Tehran
School of Electrical and Computer Engineering



Pattern Recognition

Assignment 7 (Bonus)

Due Date: 7th Bahman

Corresponding TA:

Taraneh Younesian – t.yoonesian@gmail.com

TA office hours:

Tuesdays and Wednesdays 14-16

(Room 403, Machine Learning and Computational Modeling Lab)

PROBLEM 1

Implement agglomerative hierarchical clustering using predefined python libraries (scikit-learn). The function you should use is called *AgglomerativeClustering*. Indicate confusion matrix in your report. Also compare mean of distances in each cluster and report the accuracy of your method.

PROBLEM 2

Implement sequential clustering (also known as iterative optimization) using predefined python libraries (scikit-learn). The function you should use is called *AffinityPropagation*. Report measures mentioned in problem 1.

PROBLEM 3

Calculate the within and between scatter matrices of both agglomerative hierarchical and sequential clustering. Also calculate the $trace(S_w^{-1}S_b)$ and compare the results.

PROBLEM 4

K-means is a type of optimization-based clustering in which an objective function should be minimized. It aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean, serving as a *prototype* of the cluster. Given a set of observations (x_1, x_2, \dots, x_n) where each observation is a d -dimensional real vector, k-means clustering aims to partition the n observations into k ($\leq n$) sets $S = \{S_1, S_2, \dots, S_n\}$ so as to minimize the within-cluster sum of squares (WCSS) (i.e. variance). Formally, the objective is to find:

$$\arg \min_S \sum_{i=1}^k \sum_{x \in S_i} ||x - \mu_i||^2$$

where μ_i is the mean of points in S_i :

$$\mu_i = \frac{1}{N_i} \sum_{x \in S_i} x$$

where N_i is the number of data points in S_i .

K-means algorithm:

The K -means clustering algorithm uses iterative refinement to produce a final result. The algorithm starts with initial estimates for the K centroids, which can either be randomly generated or randomly selected from the data set. The algorithm then iterates between two steps:

1. Data assignment step
2. Centroid update step

The algorithm iterates between steps one and two until a stopping criteria is met (i.e., no data points change clusters, the sum of the distances is minimized, or some maximum number of iterations is reached).

Detailed algorithm is as below:

- 1- Begin: initialize $n, k, \mu_1, \mu_2, \dots, \mu_k$
- 2- Do classify n samples according to nearest μ_i
- 3- Recompute μ_i
- 4- Until no change in μ_i
- 5- Return $\mu_1, \mu_2, \dots, \mu_k$
- 6- End

Also these links can be useful to learn this type of clustering:

<https://www.datascience.com/blog/k-means-clustering>

https://en.wikipedia.org/wiki/K-means_clustering

Implement K -means clustering (also known as iterative optimization) using predefined python libraries (which is scikit-learn and the function is called `kmeans`) and report measures mentioned in problem 1.

PROBLEM 5

Cluster validity measures describe the quality of a complete clustering. *Separation Index* is one of such measures that is proportional to the ratio of between to within distance in clusters:

$$SI = \min_j \left\{ \min_{i(i \neq j)} \left\{ \frac{d(S_i, S_j)}{\max_l d(S_l, S_l)} \right\} \right\}$$

Where

$$d(S_i, S_j) = \min\{d(x_i, x_j) \mid x_i \in S_i, x_j \in S_j\}$$

$$d(S_l, S_l) = \min\{d(x_i, x_j) \mid x_i, x_j \in S_l\}$$

Calculate this measure for problems 1, 2 and 4. Report and compare the results. Which method has better performance?