
Assignment #4

Elements of Machine Learning

Saarland University – Winter Semester 2024/25

Rabin Adhikari

7072310

raad00002@stud.uni-saarland.de

Dhimitrios Duka

7059153

dhdu00001@stud.uni-saarland.de

1 Problem 1 (K-means)

1.1 Given a dataset D of 4 points...

First, we start by calculating the distance of point x_2 and x_4 to the initial centroids \bar{x}_1 and \bar{x}_2 .

$$\begin{aligned}d_{x_2\bar{x}_1} &= \sqrt{0^2 + 3^2} = 3 \\d_{x_4\bar{x}_1} &= \sqrt{3^2 + 3^2} = 3\sqrt{2} \\d_{x_2\bar{x}_3} &= \sqrt{3^2 + 3^2} = 3\sqrt{2} \\d_{x_4\bar{x}_3} &= \sqrt{0^2 + 3^2} = 3\end{aligned}\tag{1}$$

Therefore, we can conclude that the first cluster would contain the data points x_1 and x_2 and the second cluster would contain the remaining data points x_3 and x_4 .

Now for the second iteration, first we have to calculate the new centroids.

$$\begin{aligned}\bar{x}_1 &= \left(1, \frac{5}{2}\right) \\ \bar{x}_2 &= \left(4, \frac{5}{2}\right)\end{aligned}\tag{2}$$

7 Now, we calculate the distances of every point to the new centroids.

$$\begin{aligned}
 d_{x_1 \bar{x}_1} &= \sqrt{0^2 + \left(\frac{5}{2} - 1\right)^2} = \frac{3}{2} \\
 d_{x_2 \bar{x}_1} &= \sqrt{0^2 + \left(4 - \frac{5}{2}\right)^2} = \frac{3}{2} \\
 d_{x_3 \bar{x}_1} &= \sqrt{(4 - 1)^2 + \left(1 - \frac{5}{2}\right)^2} = \frac{3}{2}\sqrt{5} \\
 d_{x_4 \bar{x}_1} &= \sqrt{(4 - 1)^2 + \left(4 - \frac{5}{2}\right)^2} = \frac{3}{2}\sqrt{5} \\
 d_{x_1 \bar{x}_2} &= \sqrt{(4 - 1)^2 + \left(\frac{5}{2} - 1\right)^2} = \frac{3}{2}\sqrt{5} \\
 d_{x_2 \bar{x}_2} &= \sqrt{(4 - 1)^2 + \left(\frac{5}{2} - 4\right)^2} = \frac{3}{2}\sqrt{5} \\
 d_{x_3 \bar{x}_2} &= \sqrt{0^2 + \left(\frac{5}{2} - 1\right)^2} = \frac{3}{2} \\
 d_{x_4 \bar{x}_2} &= \sqrt{0^2 + \left(\frac{5}{2} - 4\right)^2} = \frac{3}{2}
 \end{aligned} \tag{3}$$

8 From the above calculations, we can see that there is no reassignment of any of the datapoints to a
9 new cluster. Therefore, the algorithm has converged.

10 **1.2 Below, you are given a plot representing the within-cluster variation (also known as**
11 **inertia or within-cluster sum-of-squares, WCSS) for different numbers of clusters (k) in**
12 **k-means.**

13 **1.2.1 According to the elbow heuristic, what is the optimal number of clusters for this dataset?**
14 **Explain why did you choose this value.**

15 Based on the provided graph and the intuition behind the elbow heuristic, we would choose a value
16 of $k = 3$. The reason behind this choice is that for values smaller than 3, there is a large decrease in
17 WCSS. However, for values larger than 3, there is a much slower decrease in the WCSS, suggesting
18 diminishing returns.

19 **1.2.2 Intuitively explain how the within-cluster variation changes as the number of clusters**
20 **increases.**

21 As the number of clusters k increases, each cluster becomes smaller and more specific, thus containing
22 fewer data samples. As a result, the data samples within a cluster are closer together, therefore
23 reducing the within-cluster variation. However, based on the above exercise, we can see that the
24 within-cluster variation follows an elbow curve. This means that after a certain point, the improvement
25 of the within-cluster variation becomes smaller. This is because new clusters start splitting data points
26 that are already well-grouped together.

27 **1.2.3 Intuitively explain under what conditions the within-cluster variation equals to zero.**

28 If we suppose that $k = N$, where N is the number of data points, and we have N distinct data points,
29 the within-cluster variation would be zero since each cluster would contain only one sample and the
30 distance of that sample from itself, which is the center of the cluster, would be zero.

31 **1.2.4 The figure below shows the resulting clusters for a random dataset using both k-means**
32 **and k-medoids. Identify which of the two plots corresponds to k-medoids and explain**
33 **your reasoning.**

34 From the given plots, we can conclude that Plot 2 is the plot that corresponds to the k -medoids
35 clustering algorithm. This is due to that fact that in k -medoids, the center of the cluster is one of

36 the data samples itself, while in the k -means clustering algorithm, the center of the cluster is not
37 necessarily a sample point of the cluster.