

Artificial Intelligence and Machine Learning

Exercises – Clustering

Question 1 (KMeans algorithm) Ⓢ

Briefly describe the KMeans algorithm in your own words. Which steps have to be executed? Does the algorithm always find the optimal solution?

Question 2 (KMeans algorithm by hand) Ⓢ

Figure 1 below plots a small set of data points $x^n \in \mathbb{R}^2$ ($1 \leq n \leq 5$) as well as two randomly initialized cluster centroids μ^1 and μ^2 denoted by *. Perform one iteration of the KMeans algorithm using the Euclidean distance metric given by

$$d_2(x, \mu) := \sqrt{\sum_{m=1}^M (x_m - \mu_m)^2}.$$

Fill in the tables 1 and 2 below with your results and mark the updated cluster means in the plot. Has the algorithm already converged after one iteration?

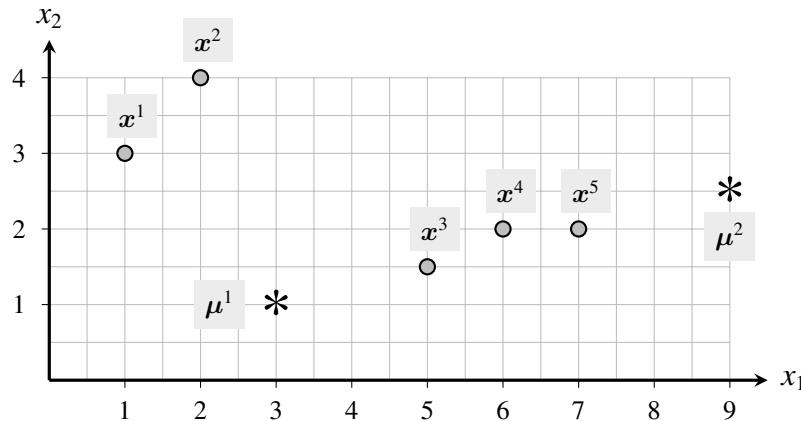


Figure 1: Plot of the dataset and the initialized cluster centroids.

n	x^n	$d_2(x^n, \mu^1)$	$d_2(x^n, \mu^2)$	Cluster assignment (1 or 2)
1	(1.0, 3.0)			
2	(2.0, 4.0)			
3	(5.0, 1.5)			
4	(6.0, 2.0)			
5	(7.0, 2.0)			

Table 1: Cluster assignments: Add your results to this table!

j	Before update: μ^j	After update: μ_{new}^j
1	(3.0, 1.0)	
2	(9.0, 2.5)	

Table 2: Update of cluster centroids: Add your results to this table!

Question 3 (Choice of K) ☒

How can you choose a suitable value for the hyperparameter K in the KMeans algorithm?



Question 4 (Implementing KMeans)

Implement the KMeans algorithm to cluster the dataset generated by the following code snippet.

```
1 from sklearn.datasets import make_blobs

3 X, _ = make_blobs(
4     n_samples=250,
5     n_features=2,
6     cluster_std=5.50,
7     centers=3,
8     shuffle=False,
9     random_state=42
10 )
```

Compress an image of your choice. The `imageio` library may be helpful.