

21/4/25

Lab-90

Date ___/___/___
Page ___

Build a K-means clustering algorithm to cluster a set of data stored in a .csv

Input

Dataset $D = \{x_1, x_2, \dots, x_n\}$ for each $x_i \in \mathbb{R}^d$, d is dimensional factor

no. of clusters: K

Output

K cluster centroids and cluster assignments for each data point

algorithm

1. Initialize centroid

• Randomly select K data points from D as initial centroids

$\mu_1, \mu_2, \dots, \mu_K$

→ choose 1st centroid from D

→ from remaining centroid

select a point t with probability

2. Iterate until convergence / max itc

• assign clusters

• for each data point x_i :

compute the euclidean distance

$$\text{dist}(x_i, \mu_j) = \sqrt{\sum_{m=1}^d (x_{i,m} - \mu_{j,m})^2}$$

2. update centroid

• for each cluster $j=1, 2, \dots, k$:

→ compute new centroid μ_j

$$\mu_j = \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$$

• check convergence

• if centroids move less than a small threshold, stop

• otherwise repeat the "assign cluster" and "update centroid" type

3. output

Return the final centroids

$\{\mu_1, \mu_2, \dots, \mu_k\}$ and cluster

assignments for each data point