

Clustering

k-means

- Random cluster centroid
- Each iteration assigning points to each centroid
- Each iteration moving the centroid to the centre of its group

$K \equiv$ number of clusters

Training set $= \{x^1, x^2, \dots, x^m\}$; $x_i \in \mathbb{R}^n$

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K$, $\mu_k \in \mathbb{R}^n$

For $i = [1, m]$

$c^i :=$ index (from 1 to K) of cluster centroid closest to $x^i \rightarrow \min_k \|x^i - \mu_k\|^2$

For $k = 1$ to K :

$\mu_k :=$ mean of points assigned to cluster k
(distortion)

• Optimization objective $\rightarrow J(c^1, \dots, c^m, \mu_1, \dots, \mu_K) = \frac{1}{2} \sum_{i=1}^m \|x^i - \mu_{c^i}\|^2$
minimize it

• $c^i \equiv$ index of cluster $[1, K]$ to which example x^i is currently assigned

• $\mu^k \equiv$ cluster centroid k , $\mu_k \in \mathbb{R}^n$

• $\mu_{c^i} \equiv$ cluster centroid of cluster to which example x^i has been assigned

- Random initialization

- should have $k < n$
- randomly pick K training examples
- set $\mu_1, \dots, \mu_k = x^{(i_1)}, \dots, x^{(i_k)}$

- choosing the value of k (the elbow)

