

K-means Clustering Algorithm

Sid

- 1) Randomly initialize 2 points : cluster centroid.
- 2) Cluster assignment : assign each data point into 2 groups based on proximity to cluster centroids.
- 3) Move centroid : compute average for all points inside each centroid group and move centroids to average.
- 4) Repeat (2-3)

Input

- K (# of clusters)
- Training Set $\{x^{(1)}, \dots, x^{(m)}\}$ $x^{(i)} \in \mathbb{R}^{(n)}$ (drop $x^{(0)}$)

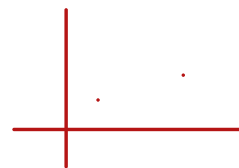
Randomly initialize K cluster centroid : μ_1, \dots, μ_K

Repeat {

Cluster assignment [for $i=1 \rightarrow m$
 $c^{(i)} := \text{index } (1-K) \text{ of cluster centroid closest to } x^{(i)}$
 $\min_k \|x^{(i)} - \mu_k\|^2$
 \downarrow
 $c^{(i)}$

Move cluster [for $k=1 \rightarrow K$
 $\mu_k := \text{avg of points assigned to cluster } k$
 $= \frac{1}{n} [x^{(k_1)} + x^{(k_2)} + \dots + x^{(k_n)}]$

}



Optimization objective

$c^{(i)}$ = index of cluster to which $x^{(i)}$ is assigned to

μ_k = cluster centroid k

$\mu_{c^{(i)}}$ = cluster centroid of cluster to which $x^{(i)}$ has been assigned

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$\min_{c, \mu} J(c, \mu)$ "distortion"
of examples

Random Initialization :

- $k < m$
- Randomly pick k training examples
- Set μ_1, \dots, μ_k = to examples

- To avoid local optima :

for $i = 1 \rightarrow 100$

randomly initialize k-means

run k-means \rightarrow compute ' c ' & ' m '

compute cost : $J(c, m)$

pick clustering w/ lowest cost.

