

Clustering method $\begin{cases} k\text{-means clustering} \\ \text{Gaussian mixture} \end{cases}$

1. k -means clustering

$D = \{X_i\}_{i=1}^N$ in \mathbb{R}^d , positive integer k , partition D into k groups.

where each data point belongs to one (and only one) cluster

one-hot vector $r_i \in \{0, 1\}^k$, $i = 1, 2, \dots, N$, $(r_i)_j = \begin{cases} 1, & \text{if } X_i \text{ is assigned to cluster } j \\ 0, & \text{otherwise} \end{cases}$

Each cluster has a "representative", $Z_j \in \mathbb{R}^d$, $j = 1, 2, \dots, k$

Determine $\{Z_j\}_{j=1}^k$ and $\{r_i\}_{i=1}^N$

The objective function is the sum of the squared distance from each data point to the center of the cluster that it is assigned to.

$$\Rightarrow \text{Min } J(\{r_i\}, \{Z_j\}) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k \frac{1}{2} (r_i)_j |X_i - Z_j|^2$$

(Note: $(r_i)_j = 1 \Leftrightarrow X_i \in \text{cluster } j$)

Step 1: Assignment step

Given $\{Z_j\}_{j=1}^k$, assign each X_i to the nearest Z_j

$$(r_i)_j = \begin{cases} 1, & \text{if } j = \arg \min_k |X_i - Z_k|^2 \\ 0, & \text{otherwise} \end{cases}$$



Step 2: Update step

Given $\{r_i\}_{i=1}^N$, set $Z_j = \frac{\sum_{i=1}^N (r_i)_j X_i}{\sum_{i=1}^N (r_i)_j}$, $j = 1, 2, \dots, k$.

$$(0 = \nabla_{Z_j} J = -\frac{1}{N} \sum_{i=1}^N (r_i)_j (X_i - Z_j) \Rightarrow 0 = \sum_{i=1}^N (r_i)_j X_i - (\sum_{i=1}^N (r_i)_j) Z_j)$$

Z_j : average (mean) of X_i 's that have been assigned to cluster j .

①: sum of samples in j -th cluster

②: number of data point in j -th cluster

Algorithm (K-means)

- Input data $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$
- Specify the number of cluster K
- Initialize the cluster centers $\{\mathbf{z}_j\}_{j=1}^K$
- Repeat:
 - update $\{\mathbf{r}_i\}_{i=1}^N$ (step 1)
 - update $\{\mathbf{z}_j\}_{j=1}^K$ (step 2)until convergence
- Return the cluster centers $\{\mathbf{z}_j\}_{j=1}^K$ and assignment $\{\mathbf{r}_i\}_{i=1}^N$

Convergence:

- For a finite dataset, the algorithm converges in a finite number of iterations.
- The algorithm may converge to a local minimum of the objective function.

Example:

