

# k-Medians algorithm

# Representative-based algorithms

- Let  $k$  be the number of clusters
- Let  $D = \{\bar{X}_1, \dots, \bar{X}_n\}$  be the dataset
- The goal is to determine  $k$  representatives  $\bar{Y}_1, \dots, \bar{Y}_k$  that minimise the following objective function

$$\sum_{i=1}^n \left[ \min_j d(\bar{X}_i, \bar{Y}_j) \right]$$

i.e. the sum of the distances of the objects to their closest representatives needs to be minimised.

# Representative-based algorithms

We obtain specific algorithms by specifying

- the way of choosing representatives, and
- the distance function  $d(\cdot, \cdot)$

In general representatives do not necessarily belong to the dataset.

# General $k$ -representatives approach

- **Initialise:** pick initial  $k$  representatives
- **Iteratively refine:**
  - **(assign step)** Assign each object to its closest representative using distance function  $d(\cdot, \cdot)$ . Denote the corresponding clusters  $C_1, \dots, C_k$
  - **(optimise step)** Determine the optimal representative  $\bar{Y}_j$  for each cluster  $C_j$  that minimises its local objective function  $\sum_{\bar{X}_i \in C_j} \left[ d(\bar{X}_i, \bar{Y}_j) \right]$

# $k$ -Means algorithm

- Representatives are chosen **not necessarily** from the dataset
- The distance function is **squared Euclidean distance**

# $k$ -Means algorithm

The objective function  $\min_{\bar{Y}_1, \dots, \bar{Y}_k} \sum_{i=1}^k \sum_{\bar{X} \in C_i} \|\bar{X} - \bar{Y}_i\|^2$

where  $C_i$  consists of the objects that are closest to  $\bar{Y}_i$ .

We want to minimise the total distance between data objects and their cluster representatives  $\bar{Y}_1, \dots, \bar{Y}_k$

This objective function is called the *within cluster sum of squares* (WCSS) objective

# $k$ -Means algorithm

Assume that the clusters  $C_1, C_2, \dots, C_k$  are fixed.

Find the set  $\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_k$  of representatives such that

$$f_{C_1, \dots, C_k}(\bar{Y}_1, \dots, \bar{Y}_k) = \sum_{i=1}^k \sum_{\bar{X} \in C_i} \|\bar{X} - \bar{Y}_i\|^2 \text{ is minimised.}$$

$$\frac{\partial f_{C_1, \dots, C_k}(\bar{Y}_1, \dots, \bar{Y}_k)}{\partial \bar{Y}_i} = - \sum_{\bar{X} \in C_i} 2(\bar{X} - \bar{Y}_i) = 0 \qquad \bar{Y}_i = \frac{1}{|C_i|} \sum_{\bar{X} \in C_i} \bar{X}$$

Just compute the centroid (mean) of each cluster and that will give you the new **optimal** cluster representatives



# $k$ -Means algorithm

What if **Manhattan distance** is more appropriate for our application than the **Euclidean distance**?

Manhattan distance (or  $L^1$  distance)

$$\text{ManDist}(\bar{X}, \bar{Y}) = \|\bar{X} - \bar{Y}\|_1 = \sum_{i=1}^d |x_i - y_i|$$



# $k$ -Medians algorithm

The objective function  $\min_{\bar{Y}_1, \dots, \bar{Y}_k} \sum_{i=1}^k \sum_{\bar{X} \in C_i} \|\bar{X} - \bar{Y}_i\|_1$

where  $C_i$  consists of the objects that are  $L^1$ -closest to  $\bar{Y}_i$ .

We want to minimise the total  $L^1$ -distance between data objects and their cluster representatives  $\bar{Y}_1, \dots, \bar{Y}_k$

# $k$ -Medians algorithm

Assume that the clusters  $C_1, C_2, \dots, C_k$  are fixed.

Find the set  $\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_k$  of representatives such that

$$f_{C_1, \dots, C_k}(\bar{Y}_1, \dots, \bar{Y}_k) = \sum_{i=1}^k \sum_{\bar{X} \in C_i} \|\bar{X} - \bar{Y}_i\|_1 \text{ is minimised.}$$

Let's consider a fixed cluster  $C$ . What is the new representative  $\bar{Y}$  that minimises  $\sum_{\bar{X} \in C} \|\bar{X} - \bar{Y}\|_1$  ?

Let  $\bar{X} = (\bar{X}^{(1)}, \bar{X}^{(2)}, \dots, \bar{X}^{(d)})$  and  $\bar{Y} = (\bar{Y}^{(1)}, \bar{Y}^{(2)}, \dots, \bar{Y}^{(d)})$ .

# $k$ -Medians algorithm

$$\arg \min_{\bar{Y}} \sum_{\bar{X} \in C} \|\bar{X} - \bar{Y}\|_1$$

$$\sum_{\bar{X} \in C} \|\bar{X} - \bar{Y}\|_1 = \sum_{\bar{X} \in C} \sum_{i=1}^d \left| \bar{X}^{(i)} - \bar{Y}^{(i)} \right| = \sum_{i=1}^d \sum_{\bar{X} \in C} \left| \bar{X}^{(i)} - \bar{Y}^{(i)} \right|$$

Since the coordinates are independent, to minimise the later sum we can minimise each of its terms independently.

More specifically, if for  $i = 1, \dots, d$ , the number  $\bar{Y}'^{(i)}$  minimises  $\sum_{\bar{X} \in C} \left| \bar{X}^{(i)} - \bar{Y}^{(i)} \right|$ , then

$$\bar{Y}' = (\bar{Y}'^{(1)}, \bar{Y}'^{(2)}, \dots, \bar{Y}'^{(d)}) \text{ minimises } \sum_{\bar{X} \in C} \|\bar{X} - \bar{Y}\|_1$$

# $k$ -Medians algorithm

Think about them as the  $i$ -th coordinates of the objects in cluster  $C$

Given  $s$  numbers  $\bar{X}_1^{(i)}, \bar{X}_2^{(i)}, \dots, \bar{X}_s^{(i)}$ , the number  
 $\bar{Y}^{(i)} = \text{median}(\bar{X}_1^{(i)}, \bar{X}_2^{(i)}, \dots, \bar{X}_s^{(i)})$

minimises  $\sum_{t=1}^s \left| \bar{X}_t^{(i)} - \bar{Y}^{(i)} \right|$ .

**Def.** A median of a sequence of numbers is any value such that at most half of the values is less than the proposed median and at most half is greater than the proposed median.

# $k$ -Medians algorithm

Hence an object  $\bar{Y}' = (\bar{Y}'^{(1)}, \bar{Y}'^{(2)}, \dots, \bar{Y}'^{(d)})$  that minimises

$$\arg \min_{\bar{Y}} \sum_{\bar{X} \in C} \|\bar{X} - \bar{Y}\|_1$$

is a **median object** of (the objects in) the cluster  $C$ ,

where  $\bar{Y}'^{(i)}$  is the median of the the  $i$ -th coordinates of the objects in the cluster  $C$ .



# $k$ -Medians algorithm

**$k$ -MediansClustering** (Number of clusters:  $k$ , Dataset:  $\{\bar{X}_1, \dots, \bar{X}_n\}$ )

## 1. Initialisation phase

Choose  $k$  cluster representatives  $\bar{Y}_1, \dots, \bar{Y}_k$  from the dataset randomly

## 2. Assignment phase

Assign all objects in the dataset to the  $L^1$ -closest representative. The resulting clusters:  $C_1, \dots, C_k$

## 3. Optimisation phase

Compute the new representatives  $\bar{Y}_1, \dots, \bar{Y}_k$  as the **medians** of the current clusters  $C_1, \dots, C_k$

Repeat Phases 2 and 3 until convergence. (Convergence is either “no objects have moved among clusters” or “fixed number of iterations specified by user”)