

# Introduction to Machine Learning

## $K$ -Means, $K$ -Meoids, $K$ -Centers and Variations

Andres Mendez-Vazquez

September 8, 2025

# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of  $K$ -Means

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of  $K$ -Means

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

# The Hardness of $K$ -means clustering

## Definition

- Given a multiset  $S \subseteq \mathbb{R}^d$ , an integer  $k$  and  $L \in \mathbb{R}$ , is there a subset  $T \subset \mathbb{R}^d$  with  $|T| = k$  such that

$$\sum_{x \in S} \min_{t \in T} \|x - t\|^2 \leq L?$$

# The Hardness of $K$ -means clustering

## Definition

- Given a multiset  $S \subseteq \mathbb{R}^d$ , an integer  $k$  and  $L \in \mathbb{R}$ , is there a subset  $T \subset \mathbb{R}^d$  with  $|T| = k$  such that

$$\sum_{x \in S} \min_{t \in T} \|x - t\|^2 \leq L?$$

## Theorem

- The  $k$ -means clustering problem is NP-complete even for  $d = 2$ .

# Reduction

## The reduction to an NP-Complete problem

- Exact Cover by 3-Sets problem

# Reduction

## The reduction to an NP-Complete problem

- Exact Cover by 3-Sets problem

## Definition

- Given a finite set  $U$  containing exactly  $3n$  elements and a collection  $\mathcal{C} = \{S_1, S_2, \dots, S_l\}$  of subsets of  $U$  each of which contains exactly 3 elements, Are there  $n$  sets in  $\mathcal{C}$  such that their union is  $U$ ?

However

There are efficient heuristic and approximation algorithms

- Which can solve this problem



# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- **$K$ -Means Clustering Heuristic**
- Convergence Criterion
- The Distance Function
- Example
- Properties of  $K$ -Means

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

# $K$ -Means - Stuart Lloyd (Circa 1957)

## History

Invented by Stuart Lloyd in Bell Labs to obtain the best quantization in a signal data set.

# $K$ -Means - Stuart Lloyd (Circa 1957)

## History

Invented by Stuart Lloyd in Bell Labs to obtain the best quantization in a signal data set.

## Something Notable

The paper was published until 1982

# K-Means - Stuart Lloyd(Circa 1957)

## History

Invented by Stuart Loyd in Bell Labs to obtain the best quantization in a signal data set.

## Something Notable

The paper was published until 1982

Basically given  $N$  vectors  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$

It tries to find  $k$  points  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k \in \mathbb{R}^d$  that minimize the expression (i.e. a partition  $S$  of the vector points):

$$\sum_{k=1}^N \sum_{i:\mathbf{x}_i \in C_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 = \sum_{k=1}^N \sum_{i:\mathbf{x}_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

# $K$ -means clustering

## $K$ -means

It is a partitional clustering algorithm.

# $K$ -means clustering

## $K$ -means

It is a partitional clustering algorithm.

### Definition

Let the set of data points (or instances)  $D$  be  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$ :

# $K$ -means clustering

## $K$ -means

It is a partitional clustering algorithm.

### Definition

Let the set of data points (or instances)  $D$  be  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$ :

- The  $K$ -means algorithm partitions the given data into  $K$  clusters.

# $K$ -means clustering

## $K$ -means

It is a partitional clustering algorithm.

### Definition

Let the set of data points (or instances)  $D$  be  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$ :

- The  $K$ -means algorithm partitions the given data into  $K$  clusters.
- Each cluster has a cluster center, called centroid.



# $K$ -means clustering

## $K$ -means

It is a partitional clustering algorithm.

### Definition

Let the set of data points (or instances)  $D$  be  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$ :

- The  $K$ -means algorithm partitions the given data into  $K$  clusters.
- Each cluster has a cluster center, called centroid.
- $K$  is specified by the user.

# $K$ -means algorithm

The  $K$ -means algorithm works as follows

Given  $k$  as the possible number of cluster:

# $K$ -means algorithm

The  $K$ -means algorithm works as follows

Given  $k$  as the possible number of cluster:

- 1 Randomly choose  $K$  data points (seeds) to be the initial **centroids**, cluster centers,
  - ▶  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

# $K$ -means algorithm

The  $K$ -means algorithm works as follows

Given  $k$  as the possible number of cluster:

- ➊ Randomly choose  $K$  data points (seeds) to be the initial **centroids**, cluster centers,
  - ▶  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$
- ➋ Assign each data point to the closest **centroid**
  - ▶  $c_i = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$

# $K$ -means algorithm

The  $K$ -means algorithm works as follows

Given  $k$  as the possible number of cluster:

- 1 Randomly choose  $K$  data points (seeds) to be the initial **centroids**, cluster centers,

- ▶  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

- 2 Assign each data point to the closest **centroid**

- ▶  $c_i = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$

- 3 Re-compute the **centroids** using the current cluster memberships.

- ▶ 
$$\mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

# $K$ -means algorithm

The  $K$ -means algorithm works as follows

Given  $k$  as the possible number of cluster:

- 1 Randomly choose  $K$  data points (seeds) to be the initial **centroids**, cluster centers,

- ▶  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

- 2 Assign each data point to the closest **centroid**

- ▶  $c_i = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$

- 3 Re-compute the **centroids** using the current cluster memberships.

- ▶ 
$$\mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

- 4 If a convergence criterion is not met, go to 2.

# What is the code trying to do?

It is trying to find a partition  $S$

$K$ -means tries to find a partition  $S$  such that it minimizes the cost function:

$$\min_S \sum_{k=1}^N \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k) \quad (1)$$

# What is the code trying to do?

It is trying to find a partition  $S$

$K$ -means tries to find a partition  $S$  such that it minimizes the cost function:

$$\min_S \sum_{k=1}^N \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k) \quad (1)$$



# What is the code trying to do?

It is trying to find a partition  $S$

$K$ -means tries to find a partition  $S$  such that it minimizes the cost function:

$$\min_S \sum_{k=1}^N \sum_{i: \mathbf{x}_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (1)$$

Where  $\boldsymbol{\mu}_k$  is the centroid for cluster  $C_k$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i: \mathbf{x}_i \in C_k} \mathbf{x}_i \quad (2)$$

Where  $N_k$  is the number of samples in the cluster  $C_k$ .

# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- **Convergence Criterion**
- The Distance Function
- Example
- Properties of  $K$ -Means

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

# What Stopping/convergence criterion should we use?

## First

No (or minimum) re-assignments of data points to different clusters.

# What Stopping/convergence criterion should we use?

## First

No (or minimum) re-assignments of data points to different clusters.

## Second

No (or minimum) change of centroids.

# What Stopping/convergence criterion should we use?

## First

No (or minimum) re-assignments of data points to different clusters.

## Second

No (or minimum) change of centroids.

## Third

Minimum decrease in the sum of squared error (SSE),

# What Stopping/convergence criterion should we use?

## First

No (or minimum) re-assignments of data points to different clusters.

## Second

No (or minimum) change of centroids.

## Third

Minimum decrease in the sum of squared error (SSE),

- $C_k$  is cluster  $k$ .

# What Stopping/convergence criterion should we use?

## First

No (or minimum) re-assignments of data points to different clusters.

## Second

No (or minimum) change of centroids.

## Third

Minimum decrease in the sum of squared error (SSE),

- $C_k$  is cluster  $k$ .
- $\mathbf{v}_k$  is the centroid of cluster  $C_k$ .

$$SSE = \sum_{k=1}^K \sum_{x \in C_k} \text{dist}(\mathbf{x}, \mathbf{v}_k)^2$$

# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- **The Distance Function**
- Example
- Properties of  $K$ -Means

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity



# The distance function

Actually, we have the following distance functions:

## Euclidean

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

# The distance function

Actually, we have the following distance functions:

## Euclidean

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

## Manhattan

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$$

# The distance function

Actually, we have the following distance functions:

## Euclidean

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

## Manhattan

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$$

## Mahalanobis

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_A = \sqrt{(\mathbf{x} - \mathbf{y})^T A (\mathbf{x} - \mathbf{y})}$$

# Outline

## 1 $K$ -Means Clustering

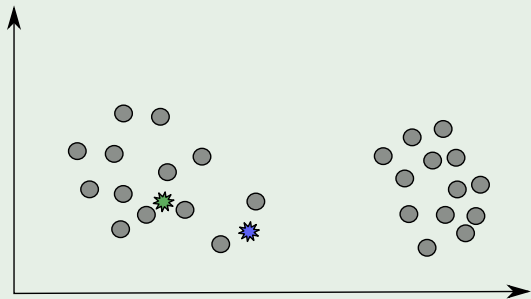
- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- **Example**
- Properties of  $K$ -Means

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

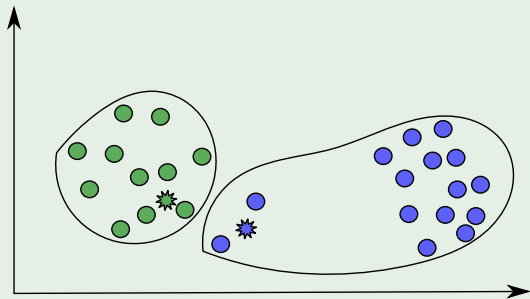
## An example

Dropping two possible centroids



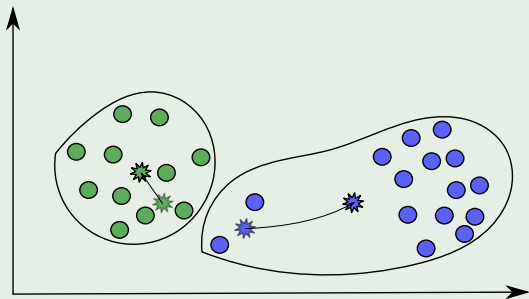
## An example

Calculate the memberships



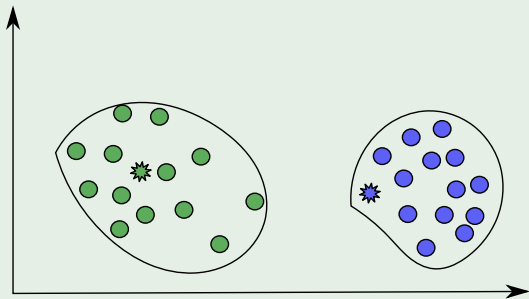
## An example

We re-calculate centroids



## An example

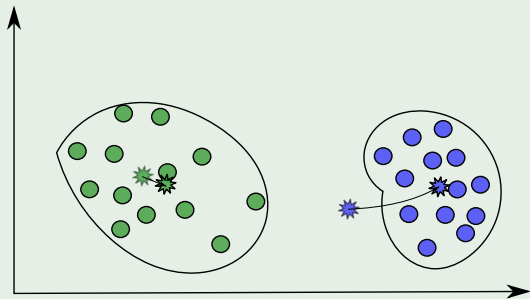
We re-calculate memberships





## An example

We re-calculate centroids and keep going



# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- **Properties of  $K$ -Means**

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

# Strengths of $K$ -means

## Strengths

- Simple: easy to understand and to implement

# Strengths of $K$ -means

## Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity:  $O(tKN)$ , where  $N$  is the number of data points,  $K$  is the number of clusters, and  $t$  is the number of iterations.

# Strengths of $K$ -means

## Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity:  $O(tKN)$ , where  $N$  is the number of data points,  $K$  is the number of clusters, and  $t$  is the number of iterations.
- Since both  $K$  and  $t$  are small.  $K$ -means is considered a linear algorithm.

# Strengths of $K$ -means

## Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity:  $O(tKN)$ , where  $N$  is the number of data points,  $K$  is the number of clusters, and  $t$  is the number of iterations.
- Since both  $K$  and  $t$  are small.  $K$ -means is considered a linear algorithm.

## Popularity

$K$ -means is the most popular clustering algorithm.

# Strengths of $K$ -means

## Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity:  $O(tKN)$ , where  $N$  is the number of data points,  $K$  is the number of clusters, and  $t$  is the number of iterations.
- Since both  $K$  and  $t$  are small.  $K$ -means is considered a linear algorithm.

## Popularity

$K$ -means is the most popular clustering algorithm.

## Note that

It terminates at a **local optimum** if SSE is used. The **global optimum** is hard to find due to complexity.

## Weaknesses of $K$ -means

### Important

The algorithm is only applicable if the mean is defined.



# Weaknesses of $K$ -means

## Important

The algorithm is only applicable if the mean is defined.

- For categorical data,  $K$ -mode - the centroid is represented by most frequent values.

# Weaknesses of $K$ -means

## Important

The algorithm is only applicable if the mean is defined.

- For categorical data,  $K$ -mode - the centroid is represented by most frequent values.

## In addition

The user needs to specify  $K$ .

# Weaknesses of $K$ -means

## Important

The algorithm is only applicable if the mean is defined.

- For categorical data,  $K$ -mode - the centroid is represented by most frequent values.

## In addition

The user needs to specify  $K$ .

## Outliers

The algorithm is sensitive to **outliers**.

# Weaknesses of $K$ -means

## Important

The algorithm is only applicable if the mean is defined.

- For categorical data,  $K$ -mode - the centroid is represented by most frequent values.

## In addition

The user needs to specify  $K$ .

## Outliers

The algorithm is sensitive to **outliers**.

- Outliers are data points that are very far away from other data points.

# Weaknesses of $K$ -means

## Important

The algorithm is only applicable if the mean is defined.

- For categorical data,  $K$ -mode - the centroid is represented by most frequent values.

## In addition

The user needs to specify  $K$ .

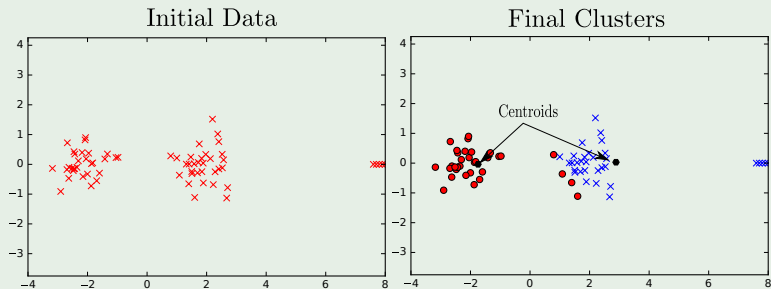
## Outliers

The algorithm is sensitive to **outliers**.

- Outliers are data points that are very far away from other data points.
- Outliers could be errors in the data recording or some special data points with very different values.

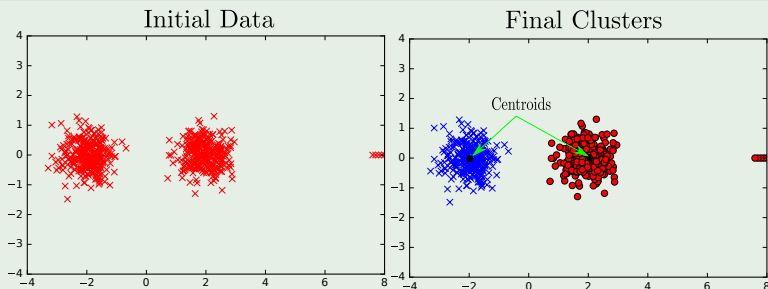
# Weaknesses of $K$ -means: Problems with outliers

## A series of outliers



# Weaknesses of $K$ -means: Problems with outliers

Nevertheless, if you have more dense clusters



## Weaknesses of $K$ -means: How to deal with outliers

### One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.



## Weaknesses of $K$ -means: How to deal with outliers

### One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

# Weaknesses of $K$ -means: How to deal with outliers

## One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

## Another method

To perform random sampling.

# Weaknesses of $K$ -means: How to deal with outliers

## One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

## Another method

To perform random sampling.

- Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.

# Weaknesses of $K$ -means: How to deal with outliers

## One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

## Another method

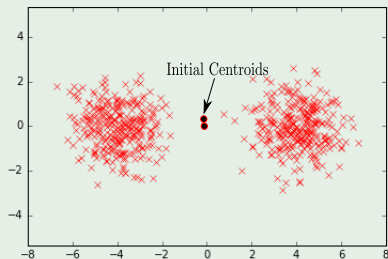
To perform random sampling.

- Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
- Assign the rest of the data points to the clusters by distance or similarity comparison, or classification.

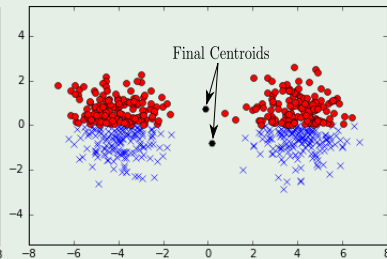
## Weaknesses of $K$ -means (cont...)

The algorithm is sensitive to **initial seeds**

Initial Centroids

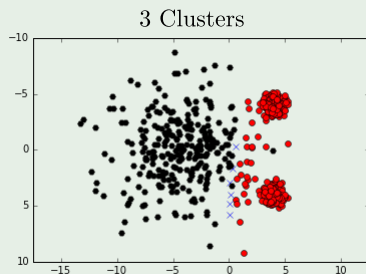
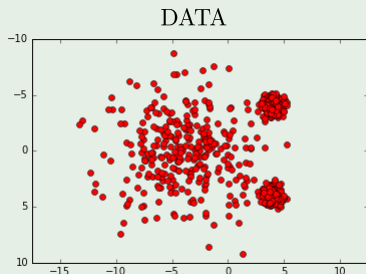


Final Clusters



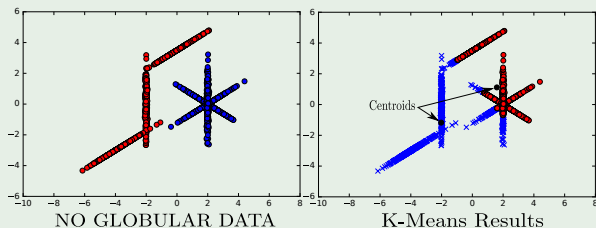
# Weaknesses of $K$ -means : Different Densities

We have three cluster nevertheless



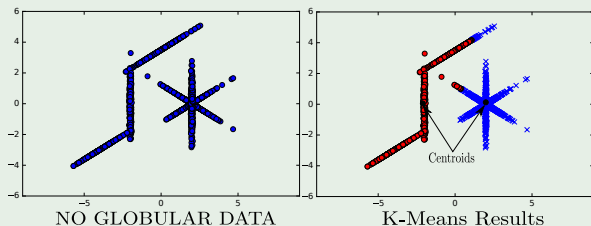
## Weaknesses of $K$ -means: Non-globular Shapes

Here, we notice that  $K$ -means may only detect globular shapes



# Weaknesses of $K$ -means: Non-globular Shapes

However, it sometimes work better than expected





# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of  $K$ -Means

## 2 $K$ -Means

- Introduction
- The Algorithm
- Complexity

Until now, we have assumed a Euclidean metric space

### Important step

- The cluster representatives  $m_1, \dots, m_k$  in are taken to be the means of the currently assigned clusters.

Until now, we have assumed a Euclidean metric space

### Important step

- The cluster representatives  $m_1, \dots, m_k$  in are taken to be the means of the currently assigned clusters.

We can generalize this by using a dissimilarity  $D(x_i, x_{i'})$

- By using an explicit optimization with respect to  $m_1, \dots, m_k$

# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of  $K$ -Means

## 2 $K$ -Means

- Introduction
- **The Algorithm**
- Complexity

# Algorithm $K$ -meoids

## Step 1

- For a given cluster assignment  $C$  find the observation in the cluster minimizing total distance to other points in that cluster:

$$i_k^* = \arg \min_{\{i | C(i)=k\}} \sum_{C(i')=k} D(\mathbf{x}_i, \mathbf{x}_{i'})$$

- ▶ Then  $m_k = \mathbf{x}_{i_k^*}$   $k = 1, \dots, K$  are the current estimates of the cluster centers.

## Step 2

- Given a current set of cluster centers  $m_1, \dots, m_k$ , minimize the total error by assigning each observation to the closest (current) cluster center:

$$C(i) = \arg \min_{1 \leq k \leq K} D(\mathbf{x}_i, m_k)$$

## Step 2

- Given a current set of cluster centers  $m_1, \dots, m_k$ , minimize the total error by assigning each observation to the closest (current) cluster center:

$$C(i) = \arg \min_{1 \leq k \leq K} D(\mathbf{x}_i, m_k)$$

## Iterate over steps 1 and 2

- Until the assignments do not change.

# Outline

## 1 $K$ -Means Clustering

- The NP-Hard Problem
- $K$ -Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of  $K$ -Means

## 2 $K$ -Means++

- Introduction
- The Algorithm
- Complexity



# Complexity

Problem, solving the first step has a complexity for  $k = 1, \dots, K$

$$O(N_k^2)$$

# Complexity

Problem, solving the first step has a complexity for  $k = 1, \dots, K$

$$O(N_k^2)$$

Given a set of cluster “centers,”  $\{i_1, i_2, \dots, i_K\}$

- Given the new assignments

$$C(i) = \arg \min_{1 \leq k \leq K} D(\mathbf{x}_i, m_k)$$

- ▶ It requires a complexity of  $O(KN)$  as before.

Therefore

We have that

- $K$ -medoids is more computationally intensive than  $K$ -means.