

# Introduction to Machine Learning

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

Andres Mendez-Vazquez

August 4, 2018

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Means

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# 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  $d = 2$ .

# 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

### 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  $U$ ?

# 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
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# 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

## Basically given $N$ vectors $x_1, \dots, x_N \in \mathbb{R}^d$

It tries to find  $k$  points  $\mu_1, \dots, \mu_k \in \mathbb{R}^d$  that minimize the expression (i.e. a partition  $S$  of the vector points):

$$\sum_{k=1}^K \sum_{i: x_i \in C_k} \|x_i - \mu_k\|^2 = \sum_{k=1}^K \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k)$$

# 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

Basically given  $N$  vectors  $x_1, \dots, x_N \in \mathbb{R}^d$

It tries to find  $k$  points  $\mu_1, \dots, \mu_k \in \mathbb{R}^d$  that minimize the expression (i.e. a partition  $S$  of the vector points):

$$\sum_{k=1}^K \sum_{i: x_i \in C_k} \|x_i - \mu_k\|^2 = \sum_{k=1}^K \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k)$$

# 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

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  $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$ :

- 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 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 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 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:

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

- ▶  $\{v_1, \dots, v_k\}$

- 2 Assign each data point to the closest centroid

- ▶  $c_i = \arg \min_j \{dist(x_i - v_j)\}$

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

- ▶ 
$$v_j = \frac{\sum_{i=1}^n I(c_i = j) x_i}{\sum_{i=1}^n I(c_i = j)}$$

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

# $K$ -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.

# $K$ -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.

# $K$ -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.

# $K$ -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:

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

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  $C_k$ .

## 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  $C_k$ .



## 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
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# 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),

- $C_k$  is cluster  $k$ .
- $v_k$  is the centroid of cluster  $C_k$ .

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

# 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  $k$ .

- $v_k$  is the centroid of cluster  $C_k$ .

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

# 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
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function



# 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

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

# 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

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

# 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
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

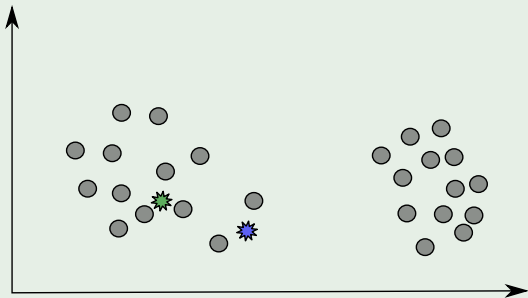
- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

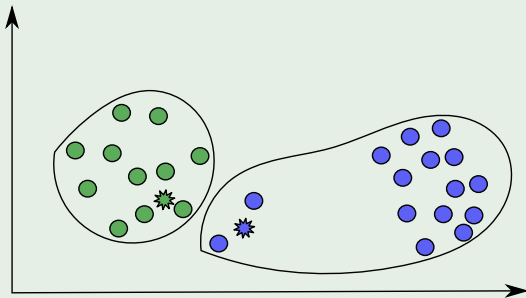
## An example

Dropping two possible centroids



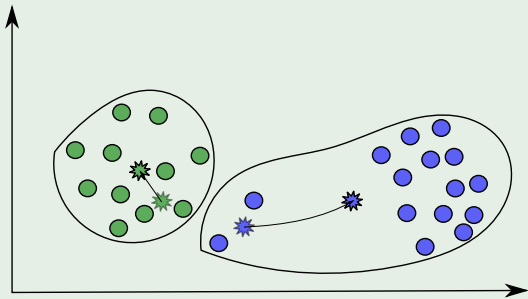
# An example

Calculate the memberships



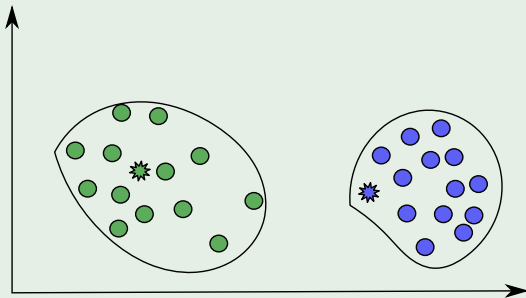
## An example

We re-calculate centroids



## An example

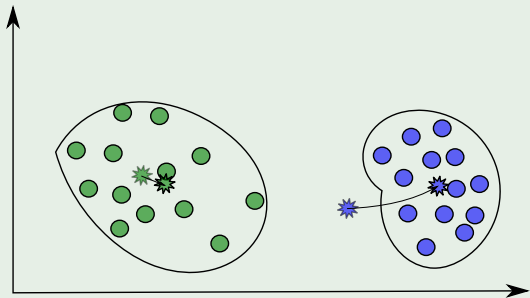
We re-calculate memberships





## 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**
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# 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.

## Weakness

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

# 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.

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

# 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.

- 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.

- 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

## 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

## 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

## 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  $K$ .

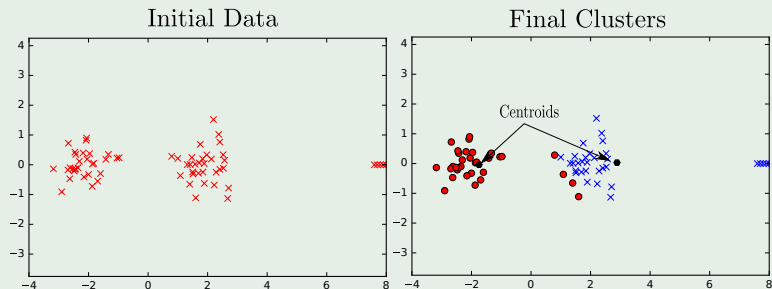
## Outliers

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.

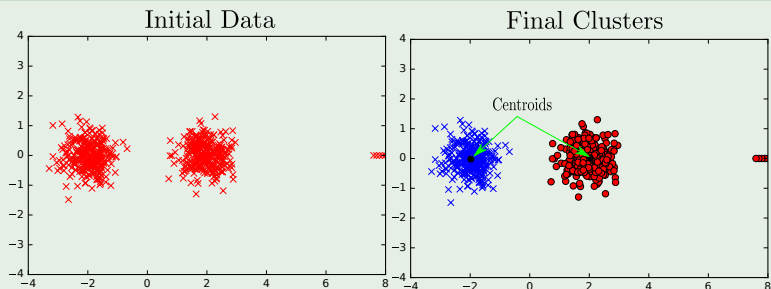
# Weaknesses of $K$ -means: Problems with outliers

## 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.

- 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.

- 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: 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: 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: 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

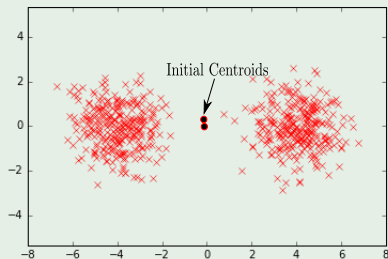
To perform random sampling.

- 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.

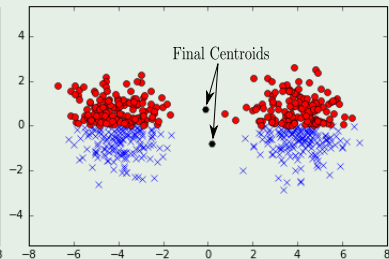
## Weaknesses of $K$ -means (cont...)

The algorithm is sensitive to **initial seeds**

Initial Centroids

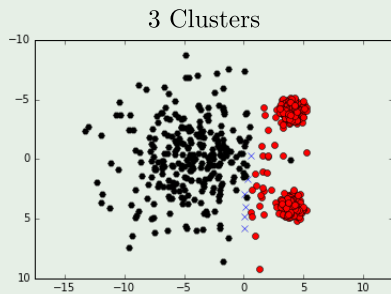
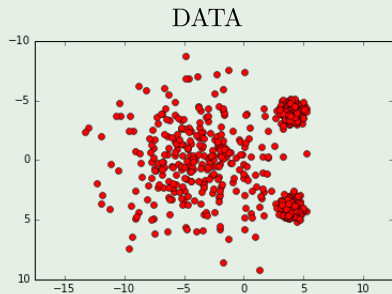


Final Clusters



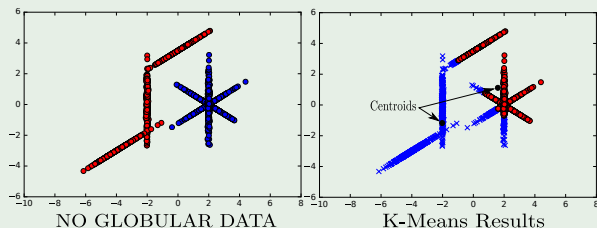
## Weaknesses of $K$ -means : Different Densities

We have three cluster nevertheless



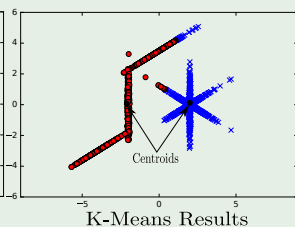
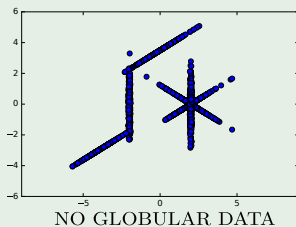
## Weaknesses of $K$ -means: Non-globular Shapes

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
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

## Consider the following

### Theorem

- Every matrix  $A \in R^{m \times n}$  has an SVD.

### Frobenius Matrix Norm

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2} = \sqrt{\text{trace}(A^T A)}$$

## Consider the following

### Theorem

- Every matrix  $A \in R^{m \times n}$  has an SVD.

### Frobenius Matrix Norm

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2} = \sqrt{\text{trace}(A^T A)}$$

## Then, you have a the Eckhart-Young Theorem

### Theorem

- Let  $A$  be a real  $m \times n$  matrix. Then for any  $k \in \mathbb{N}$  and any  $m \times m$  orthogonal projection matrix  $P$  of rank  $k$ , we have

$$\|A - P_k A\|_F \leq \|A - PA\|_F$$

- ▶ with  $P_k = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^T$

Thus

We have the Covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Therefore, we have the following decomposition

$$S = U \Sigma U^T$$

- Where  $UU^T = I$  and  $U$  is a  $d \times d$  matrix

Thus

We have the Covariance matrix

$$S = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Therefore, we have the following decomposition

$$S = U\Sigma U^T$$

- Where  $UU^T = I$  and  $U$  is a  $d \times d$  matrix

# Orthogonal Projection

Therefore, we have that  $U$  is an orthogonal projection

- Given that  $UU^T = I$  and  $Ux = x$

Now, we can rewrite  $f_{k\text{-means}}$

$$f_{k\text{-means}} = \min_{\mu_1, \dots, \mu_k} \sum_{i \in [n]} \min_{j \in [k]} \|x_i - \mu_j\|^2$$

# Orthogonal Projection

Therefore, we have that  $U$  is a orthogonal projection

- Given that  $UU^T = I$  and  $U\mathbf{x} = \mathbf{x}$

Now, we can re-write  $k$ -means

$$f_{k\text{-mean}} = \min_{\mu_1, \dots, \mu_k} \sum_{i \in [n]} \min_{j \in [k]} \|\mathbf{x}_i - \mu_j\|^2$$



Then

PCA can also re-write the cost function

$$f_{PCA} = \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 = \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2$$

Where

- Given that  $P_k$  is a projection into dimension  $k$  and  $\mathbf{y} \in P_k$  means that  $P_k \mathbf{y} = \mathbf{y}$

Furthermore

$$\arg \min_{\mathbf{y} \in P} \|\mathbf{x} - \mathbf{y}\| = P\mathbf{x}$$

Then

PCA can also re-write the cost function

$$f_{PCA} = \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 = \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2$$

Where

- Given that  $P_k$  is a projection into dimension  $k$  and  $\mathbf{y} \in P_k$  means that  $P_k \mathbf{y} = \mathbf{y}$

Furthermore

$$\arg \min_{\mathbf{y} \in P} \|\mathbf{x} - \mathbf{y}\| = P\mathbf{x}$$

Then

PCA can also re-write the cost function

$$f_{PCA} = \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 = \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2$$

Where

- Given that  $P_k$  is a projection into dimension  $k$  and  $\mathbf{y} \in P_k$  means that  $P_k \mathbf{y} = \mathbf{y}$

Furthermore

$$\arg \min_{\mathbf{y} \in P} \|\mathbf{x} - \mathbf{y}\| = P\mathbf{x}$$

## Thus, using the Eckhart-Young Theorem

Assume  $P_k^*$  which contains the  $k$  optimal centers

- Given that  $\mu_j \in P_k^*$

$$\begin{aligned} f_{k\text{-mean}} &= \sum_{i \in [n]} \min_{j \in [k]} \|\mathbf{x}_i - \mu_j^*\|^2 \\ &\geq \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k^*} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &\geq \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &= \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 \\ &= f_{PCA} \end{aligned}$$

## Thus, using the Eckhart-Young Theorem

Assume  $P_k^*$  which contains the  $k$  optimal centers

- Given that  $\mu_j \in P_k^*$

$$\begin{aligned} f_{k\text{-mean}} &= \sum_{i \in [n]} \min_{j \in [k]} \|\mathbf{x}_i - \mu_j^*\|^2 \\ &\geq \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k^*} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &\geq \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &= \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 \\ &= f_{PCA} \end{aligned}$$

## Thus, using the Eckhart-Young Theorem

Assume  $P_k^*$  which contains the  $k$  optimal centers

- Given that  $\mu_j \in P_k^*$

$$\begin{aligned} f_{k\text{-mean}} &= \sum_{i \in [n]} \min_{j \in [k]} \|\mathbf{x}_i - \mu_j^*\|^2 \\ &\geq \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k^*} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &\geq \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &= \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 \\ &= f_{POA} \end{aligned}$$

## Thus, using the Eckhart-Young Theorem

Assume  $P_k^*$  which contains the  $k$  optimal centers

- Given that  $\mu_j \in P_k^*$

$$\begin{aligned} f_{k\text{-mean}} &= \sum_{i \in [n]} \min_{j \in [k]} \|\mathbf{x}_i - \mu_j^*\|^2 \\ &\geq \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k^*} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &\geq \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &= \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 \end{aligned}$$

= FPOA

## Thus, using the Eckhart-Young Theorem

Assume  $P_k^*$  which contains the  $k$  optimal centers

- Given that  $\mu_j \in P_k^*$

$$\begin{aligned} f_{k\text{-mean}} &= \sum_{i \in [n]} \min_{j \in [k]} \|\mathbf{x}_i - \mu_j^*\|^2 \\ &\geq \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k^*} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &\geq \min_{P_k} \sum_{i \in [n]} \min_{\mathbf{y}_i \in P_k} \|\mathbf{x}_i - \mathbf{y}_i\|^2 \\ &= \min_{P_k} \sum_{i \in [n]} \|\mathbf{x}_i - P_k \mathbf{x}_i\|^2 \\ &= f_{PCA} \end{aligned}$$



## Therefore

Now, consider solving  $k$ -means on the points  $\mathbf{y}_i$  instead

- They are embedded into dimension exactly  $k$  by projection  $P_k$

Therefore, given  $\mathcal{X} = \{\mathbf{x}_i\}$  and  $\mathcal{Y} = \{\mathbf{y}_i\}$ ,

- Where the  $\hat{S}$  and  $\hat{\mu}$  are the assignments and centers of the projected points  $\mathbf{y}_i$ :

# Therefore

Now, consider solving  $k$ -means on the points  $\mathbf{y}_i$  instead

- They are embedded into dimension exactly  $k$  by projection  $P_k$

Therefore, given  $P\mathbf{x}_i = \mathbf{y}_i$  and  $\hat{\mu}_j = P\mu_j$

- Where the  $\hat{S}$  and  $\hat{\mu}$  are the assignments and centers of the projected points  $\mathbf{y}_i$ :

$$\sum_{j \in [k]} \sum_{i \in S_j} \|\mathbf{x}_i - \mu_j\|^2 \geq \sum_{j \in [k]} \sum_{i \in S_j} \|P\mathbf{x}_i - P\mu_j\|^2$$

$$= \sum_{j \in [k]} \sum_{i \in S_j} \|\mathbf{y}_i - \hat{\mu}_j\|^2$$

$$\geq \sum_{j \in [k]} \sum_{i \in \hat{S}_j} \|\mathbf{y}_i - \hat{\mu}_j\|^2 = f_k^* \text{-means}$$

# Therefore

Now, consider solving  $k$ -means on the points  $\mathbf{y}_i$  instead

- They are embedded into dimension exactly  $k$  by projection  $P_k$

Therefore, given  $P\mathbf{x}_i = \mathbf{y}_i$  and  $\hat{\mu}_j = P\mu_j$

- Where the  $\hat{S}$  and  $\hat{\mu}$  are the assignments and centers of the projected points  $\mathbf{y}_i$ :

$$\begin{aligned} \sum_{j \in [k]} \sum_{i \in S_j} \|\mathbf{x}_i - \mu_j\|^2 &\geq \sum_{j \in [k]} \sum_{i \in S_j} \|P\mathbf{x}_i - P\mu_j\|^2 \\ &= \sum_{j \in [k]} \sum_{i \in S_j} \|\mathbf{y}_i - \hat{\mu}_j\|^2 \end{aligned}$$

$$\geq \sum_{j \in [k]} \sum_{i \in \hat{S}_j} \|\mathbf{y}_i - \hat{\mu}_j\|^2 = f_k^* \text{-means}$$

# Therefore

Now, consider solving  $k$ -means on the points  $\mathbf{y}_i$  instead

- They are embedded into dimension exactly  $k$  by projection  $P_k$

Therefore, given  $P\mathbf{x}_i = \mathbf{y}_i$  and  $\hat{\mu}_j = P\mu_j$

- Where the  $\hat{S}$  and  $\hat{\mu}$  are the assignments and centers of the projected points  $\mathbf{y}_i$ :

$$\begin{aligned} \sum_{j \in [k]} \sum_{i \in S_j} \|\mathbf{x}_i - \mu_j\|^2 &\geq \sum_{j \in [k]} \sum_{i \in S_j} \|P\mathbf{x}_i - P\mu_j\|^2 \\ &= \sum_{j \in [k]} \sum_{i \in S_j} \|\mathbf{y}_i - \hat{\mu}_j\|^2 \\ &\geq \sum_{j \in [k]} \sum_{i \in \hat{S}_j} \|\mathbf{y}_i - \hat{\mu}_j\|^2 = f_{k\text{-means}}^* \end{aligned}$$

Therefore, your best bet

## Steps

- 1 Compute the PCA of the points  $x_i$  into dimension  $k$ .
- 2 Solve  $k$ -means on the points  $y_i$  in dimension  $k$ .
- 3 Output the resulting clusters and centers.

Therefore, your best bet

## Steps

- 1 Compute the PCA of the points  $x_i$  into dimension  $k$ .
- 2 Solve  $k$ -means on the points  $y_i$  in dimension  $k$ .
- 3 Output the resulting clusters and centers.

Therefore, your best bet

## Steps

- 1 Compute the PCA of the points  $x_i$  into dimension  $k$ .
- 2 Solve  $k$ -means on the points  $y_i$  in dimension  $k$ .
- 3 Output the resulting clusters and centers.

Given that

We have that

$$f_{new} = \sum_{j \in [k]} \sum_{i \in S_j^*} \left\| \mathbf{x}_i - \mu_j^* \right\|^2 = *$$

Therefore by the fact that  $\mathbf{x}_i - \mathbf{y}_i$  and  $\mathbf{y}_i - \mu_j^*$  are perpendicular

$$* = \sum_{j \in [k]} \sum_{i \in S_j^*} \left\{ \left\| \mathbf{x}_i - \mathbf{y}_i \right\|^2 + \left\| \mathbf{y}_i - \mu_j^* \right\|^2 \right\} = **$$

Finally

$$** = \sum_{i \in [n]} \left\| \mathbf{x}_i - \mathbf{y}_i \right\|^2 + \sum_{j \in [k]} \sum_{i \in S_j^*} \left\| \mathbf{y}_i - \mu_j^* \right\|^2$$



Given that

We have that

$$f_{new} = \sum_{j \in [k]} \sum_{i \in S_j^*} \left\| \mathbf{x}_i - \mu_j^* \right\|^2 = *$$

Therefore by the fact that  $\mathbf{x}_i - \mathbf{y}_i$  and  $\mathbf{y}_i - \mu_j^*$  are perpendiculars

$$* = \sum_{j \in [k]} \sum_{i \in S_j^*} \left\{ \left\| \mathbf{x}_i - \mathbf{y}_i \right\|^2 + \left\| \mathbf{y}_i - \mu_j^* \right\|^2 \right\} = **$$

Finally

$$** = \sum_{i \in [n]} \left\| \mathbf{x}_i - \mathbf{y}_i \right\|^2 + \sum_{j \in [k]} \sum_{i \in S_j^*} \left\| \mathbf{y}_i - \mu_j^* \right\|^2$$

Given that

We have that

$$f_{new} = \sum_{j \in [k]} \sum_{i \in S_j^*} \left\| \mathbf{x}_i - \mu_j^* \right\|^2 = *$$

Therefore by the fact that  $\mathbf{x}_i - \mathbf{y}_i$  and  $\mathbf{y}_i - \mu_j^*$  are perpendiculars

$$* = \sum_{j \in [k]} \sum_{i \in S_j^*} \left\{ \left\| \mathbf{x}_i - \mathbf{y}_i \right\|^2 + \left\| \mathbf{y}_i - \mu_j^* \right\|^2 \right\} = **$$

Finally

$$** = \sum_{i \in [n]} \left\| \mathbf{x}_i - \mathbf{y}_i \right\|^2 + \sum_{j \in [k]} \sum_{i \in S_j^*} \left\| \mathbf{y}_i - \mu_j^* \right\|^2$$

Therefore, we have

### Something Notable

$$f_{PCA} + f_{k\text{-means}}^* \leq 2f_{k\text{-means}}$$

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- **Introduction**
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

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  $\phi(x, y)$

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

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  $D(\mathbf{x}_i, \mathbf{x}_{i'})$

- 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
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Meoids

- Introduction
- **The Algorithm**
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# 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:

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

Iterate over steps 1 and 2

- Until the assignments do not change.

## 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
- $K$ -Means and Principal Component Analysis

## 2 $K$ -Means

- Introduction
- The Algorithm
- **Complexity**

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# Complexity

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

$$O(N_k^2)$$

Given a set of cluster centers  $\{m_1, \dots, m_K\}$

- Given the new assignments

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

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

# 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.

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- **Introduction**
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# The $K$ -center Problem

## The input

It is a set of points with distances represented by a weighted graph  $G = (V, V \times V)$ .



# The $K$ -center Problem

## The input

It is a set of points with distances represented by a weighted graph  $G = (V, V \times V)$ .

## Output

- Select  $K$  centroids in  $G$ .
- Such that minimize maximum distance of every point to a centroid.

# The $K$ -center Problem

## The input

It is a set of points with distances represented by a weighted graph  $G = (V, V \times V)$ .

## Output

- Select  $K$  centroids in  $G$ .
- Such that minimize maximum distance of every point to a centroid.

Theorem (In the general case for any distance)

It is NP-hard to approximate the general  $K$ -center problem within any factor  $\alpha$ .

# The $K$ -center Problem

## The input

It is a set of points with distances represented by a weighted graph  $G = (V, V \times V)$ .

## Output

- Select  $K$  centroids in  $G$ .
- Such that minimize maximum distance of every point to a centroid.

## Theorem (In the general case for any distance)

It is NP-hard to approximate the general  $K$ -center problem within any factor  $\alpha$ .

Therefore

We change the distance to be constrained by

The Triangle Inequality

Given  $x, y$  and  $z$

$$L(x, z) \leq L(x, y) + L(y, z)$$

Therefore

We change the distance to be constrained by

The Triangle Inequality

Given  $x, y$  and  $z$

$$L(x, z) \leq L(x, y) + L(y, z)$$

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- **Re-Stating the  $K$ -center as a Clustering Problem**
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

## We have a new criterion

Instead of using the  $K$ -mean criterion

We use the  $K$ -center criterion under the triangle inequality.

### New Criterion

The  $K$ -center criterion partitions the points into  $K$  clusters so as to minimize the maximum distance of any point to its cluster center.

## We have a new criterion

Instead of using the  $K$ -mean criterion

We use the  $K$ -center criterion under the triangle inequality.

### New Criterion

The  $K$ -center criterion partitions the points into  $K$  clusters so as to minimize the maximum distance of any point to its cluster center.



# Explanation

## Setup

Suppose we have a data set  $A$  that contains  $N$  objects.

# Explanation

## Setup

Suppose we have a data set  $A$  that contains  $N$  objects.

## We want the following

We want to partition  $A$  into  $K$  sets labeled  $C_1, C_2, \dots, C_K$ .

# Explanation

## Setup

Suppose we have a data set  $A$  that contains  $N$  objects.

## We want the following

We want to partition  $A$  into  $K$  sets labeled  $C_1, C_2, \dots, C_K$ .

## Now

Define a cluster size for any cluster  $C_k$  as follows.

- The smallest value  $D$  for which all points in this cluster  $C_k$  are:
  - ▶ Within distance  $D$  of each other.
  - ▶ Or within distance  $\frac{D}{2}$  of some point called the cluster center.

# Explanation

## Setup

Suppose we have a data set  $A$  that contains  $N$  objects.

## We want the following

We want to partition  $A$  into  $K$  sets labeled  $C_1, C_2, \dots, C_K$ .

## Now

Define a cluster size for any cluster  $C_k$  as follows.

- The smallest value  $D$  for which all points in this cluster  $C_k$  are:
  - ▶ Within distance  $D$  of each other.
  - ▶ Or within distance  $\frac{D}{2}$  of some point called the cluster center.

# Explanation

## Setup

Suppose we have a data set  $A$  that contains  $N$  objects.

## We want the following

We want to partition  $A$  into  $K$  sets labeled  $C_1, C_2, \dots, C_K$ .

## Now

Define a cluster size for any cluster  $C_k$  as follows.

- The smallest value  $D$  for which all points in this cluster  $C_k$  are:
  - ▶ Within distance  $D$  of each other.
  - ▶ Or within distance  $\frac{D}{2}$  of some point called the cluster center.

# Explanation

## Setup

Suppose we have a data set  $A$  that contains  $N$  objects.

## We want the following

We want to partition  $A$  into  $K$  sets labeled  $C_1, C_2, \dots, C_K$ .

## Now

Define a cluster size for any cluster  $C_k$  as follows.

- The smallest value  $D$  for which all points in this cluster  $C_k$  are:
  - ▶ Within distance  $D$  of each other.
  - ▶ Or within distance  $\frac{D}{2}$  of some point called the cluster center.

## Another Way

### We have

Another way to define the cluster size:

- $D$  is the maximum pairwise distance between an arbitrary pair of points in the cluster.
- Or twice the maximum distance between a data point and a chosen centroid.

### Thus

Denoting the cluster size of  $C_k$  by  $D_k$ , we have that the cluster size of partition (the way the points are grouped)  $S$  by:

$$D = \max_{k=1, \dots, K} D_k \quad (3)$$

### In other words

The cluster size of a partition composed of multiple clusters is the maximum size of these clusters.

## Another Way

### We have

Another way to define the cluster size:

- $D$  is the maximum pairwise distance between an arbitrary pair of points in the cluster.
- Or twice the maximum distance between a data point and a chosen centroid.

### Thus

Denoting the cluster size of  $C_k$  by  $D_k$ , we have that the cluster size of partition (the way the points are grouped)  $S$  by:

$$D = \max_{k=1, \dots, K} D_k \quad (3)$$

### In another words

The cluster size of a partition composed of multiple clusters is the maximum size of these clusters.



## Another Way

### We have

Another way to define the cluster size:

- $D$  is the maximum pairwise distance between an arbitrary pair of points in the cluster.
- Or twice the maximum distance between a data point and a chosen centroid.

### Thus

Denoting the cluster size of  $C_k$  by  $D_k$ , we have that the cluster size of partition (the way the points are grouped)  $S$  by:

$$D = \max_{k=1,\dots,K} D_k \quad (3)$$

### In another words

The cluster size of a partition composed of multiple clusters is the maximum size of these clusters.

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- **Comparison with  $K$ -means**
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

## Comparison with $K$ -means

We use the following distance for comparison

In order to compare these methods, we use Euclidean distance.

In  $K$ -means we assume the distance between vectors is the squared Euclidean distance

$K$ -means tries to find a partition  $S$  that minimizes:

$$\min_S \sum_{k=1}^K \sum_{i: x_i \in C_k} (\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k) \quad (4)$$

where  $\mu_k$  is the centroid for cluster

$$\mu_k = \frac{1}{N_k} \sum_{i: x_i \in C_k} \mathbf{x}_i \quad (5)$$

## Comparison with $K$ -means

We use the following distance for comparison

In order to compare these methods, we use Euclidean distance.

In  $K$ -means we assume the distance between vectors is the squared Euclidean distance

$K$ -means tries to find a partition  $S$  that minimizes:

$$\min_S \sum_{k=1}^N \sum_{i: x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (4)$$

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

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i: x_i \in C_k} \mathbf{x}_i \quad (5)$$

## Comparison with $K$ -means

We use the following distance for comparison

In order to compare these methods, we use Euclidean distance.

In  $K$ -means we assume the distance between vectors is the squared Euclidean distance

$K$ -means tries to find a partition  $S$  that minimizes:

$$\min_S \sum_{k=1}^N \sum_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (4)$$

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

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i:x_i \in C_k} \mathbf{x}_i \quad (5)$$

Now, what about the  $K$ -centers

$K$ -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

## Now, what about the $K$ -centers

$K$ -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1, \dots, K} \max_{i: x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

Where

$\boldsymbol{\mu}_k$  is called the “centroid”, but may not be the mean vector.

## Now, what about the $K$ -centers

$K$ -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

### Where

$\boldsymbol{\mu}_k$  is called the “centroid”, but may not be the mean vector.

### Properties

- The above objective function shows that for each cluster, only the worst scenario matters, that is, the farthest data point to the centroid.
- Moreover, among the clusters, only the worst cluster matters, whose farthest data point yields the maximum distance to the centroid comparing with the farthest data points of the other clusters.



## Now, what about the $K$ -centers

$K$ -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

### Where

$\boldsymbol{\mu}_k$  is called the “centroid”, but may not be the mean vector.

### Properties

- The above objective function shows that for each cluster, only the worst scenario matters, that is, the farthest data point to the centroid.
- Moreover, among the clusters, only the worst cluster matters, whose farthest data point yields the maximum distance to the centroid comparing with the farthest data points of the other clusters.

## In other words

### First

This minimax type of problem is much harder to solve than solving the objective function of k-means.

Another formulation of  $k$ -center minimizes the worst case pairwise distance instead of using centroids.

$$\min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (7)$$

### With

$L(x_i, x_j)$  denotes any distance between a pair of objects in the same cluster.

## In other words

### First

This minimax type of problem is much harder to solve than solving the objective function of k-means.

Another formulation of  $K$ -center minimizes the worst case pairwise distance instead of using centroids

$$\min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (7)$$

$L(x_i, x_j)$  denotes any distance between a pair of objects in the same cluster.

## In other words

### First

This minimax type of problem is much harder to solve than solving the objective function of k-means.

Another formulation of  $K$ -center minimizes the worst case pairwise distance instead of using centroids

$$\min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (7)$$

### With

$L(x_i, x_j)$  denotes any distance between a pair of objects in the same cluster.

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- **The Greedy  $K$ -Center Algorithm**
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# Greedy Algorithm for $K$ -Center

## Main Idea

The idea behind the Greedy Algorithm is to choose a subset  $H$  from the original dataset  $S$  consisting of  $K$  points that are farthest apart from each other.

## Intuition

Since the points in set  $H$  are far apart then the worst-case scenario has been taken care of and hopefully the cluster size for the partition is small.

## What

Each point  $h_k \in H$  represents one cluster or subset of points  $C_k$ .

# Greedy Algorithm for $K$ -Center

## Main Idea

The idea behind the Greedy Algorithm is to choose a subset  $H$  from the original dataset  $S$  consisting of  $K$  points that are farthest apart from each other.

## Intuition

Since the points in set  $H$  are far apart then the worst-case scenario has been taken care of and hopefully the cluster size for the partition is small.

Each point  $h_k \in H$  represents one cluster or subset of points  $C_k$ .

# Greedy Algorithm for $K$ -Center

## Main Idea

The idea behind the Greedy Algorithm is to choose a subset  $H$  from the original dataset  $S$  consisting of  $K$  points that are farthest apart from each other.

## Intuition

Since the points in set  $H$  are far apart then the worst-case scenario has been taken care of and hopefully the cluster size for the partition is small.

## Thus

Each point  $h_k \in H$  represents one cluster or subset of points  $C_k$ .



# Then

## Something Notable

We can think of it as a centroid.

### However

Technically it is not a centroid because it tends to be at the boundary of a cluster, but conceptually we can think of it as a centroid.

### Important

The way that we partition these points given the centroids is the same as in  $K$ -means, that is, the nearest-neighbor rule.

Then

## Something Notable

We can think of it as a centroid.

## However

Technically it is not a centroid because it tends to be at the boundary of a cluster, but conceptually we can think of it as a centroid.

## Implication

The way that we partition these points given the centroids is the same as in  $K$ -means, that is, the nearest-neighbor rule.

# Then

## Something Notable

We can think of it as a centroid.

## However

Technically it is not a centroid because it tends to be at the boundary of a cluster, but conceptually we can think of it as a centroid.

## Important

The way that we partition these points given the centroids is the same as in  $K$ -means, that is, the nearest-neighbor rule.

## Specifically

### We do the following

For every point  $x_i$ , in order to see which cluster  $C_k$  it is partitioned into, we compute its distance to each cluster centroid as follows, and find out which centroid is the closest:

$$L(x_i, h_k) = \min_{k'=1, \dots, K} L(x_i, h_{k'}) \quad (8)$$

### Thus

Whichever centroid with the minimum distance is selected as the cluster for  $x_i$ .

## Specifically

### We do the following

For every point  $x_i$ , in order to see which cluster  $C_k$  it is partitioned into, we compute its distance to each cluster centroid as follows, and find out which centroid is the closest:

$$L(x_i, h_k) = \min_{k'=1, \dots, K} L(x_i, h_{k'}) \quad (8)$$

### Thus

Whichever centroid with the minimum distance is selected as the cluster for  $x_i$ .

# Important

## For K-center clustering

We only need pairwise distance  $L(\mathbf{x}_i, \mathbf{x}_j)$  for any  $\mathbf{x}_i, \mathbf{x}_j \in S$ .

Where

$\mathbf{x}_j$  can be a non-vector representation of the objects.

As long as we can calculate

$L(\mathbf{x}_i, \mathbf{x}_j)$  which makes the  $K$ -center more general than the  $K$ -means.

# Important

## For K-center clustering

We only need pairwise distance  $L(\mathbf{x}_i, \mathbf{x}_j)$  for any  $\mathbf{x}_i, \mathbf{x}_j \in S$ .

## Where

$\mathbf{x}_i$  can be a non-vector representation of the objects.

As long as we can calculate

$L(\mathbf{x}_i, \mathbf{x}_j)$  which makes the  $K$ -center more general than the  $K$ -means.

# Important

## For $K$ -center clustering

We only need pairwise distance  $L(\mathbf{x}_i, \mathbf{x}_j)$  for any  $\mathbf{x}_i, \mathbf{x}_j \in S$ .

## Where

$\mathbf{x}_i$  can be a non-vector representation of the objects.

## As long we can calculate

$L(\mathbf{x}_i, \mathbf{x}_j)$  which makes the  $K$ -center more general than the  $K$ -means.



# Properties

## Something Notable

The greedy algorithm achieves an approximation factor of 2 as the distance measure  $L$  satisfies the triangle inequality.

Thus, we have that

$$D^* = \min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (9)$$

Then, we have the following guarantee for the greedy algorithm

$$D \leq 2D^* \quad (10)$$

# Properties

## Something Notable

The greedy algorithm achieves an approximation factor of 2 as the distance measure  $L$  satisfies the triangle inequality.

Thus, we have that

$$D^* = \min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (9)$$

Then, we have the following guarantee for the greedy algorithm:

$$D \leq 2D^* \quad (10)$$

# Properties

## Something Notable

The greedy algorithm achieves an approximation factor of 2 as the distance measure  $L$  satisfies the triangle inequality.

Thus, we have that

$$D^* = \min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (9)$$

Then, we have the following guarantee for the greedy algorithm

$$D \leq 2D^* \quad (10)$$

# Nevertheless

We have that

$K$ -center does not provide a locally optimal solution.

We can get only a solution

Guaranteeing to be within a certain performance range of the theoretical optimal solution.

# Nevertheless

We have that

$K$ -center does not provide a locally optimal solution.

We can get only a solution

Guaranteeing to be within a certain performance range of the theoretical optimal solution.

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- **Pseudo-Code**
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# Setup

## First

Set  $H$  denotes the set of cluster centroids or cluster of representative objects  $\{\mathbf{h}_1, \dots, \mathbf{h}_k\} \subset S$ .

## Second

Let *cluster* ( $x_i$ ) be the identity of the cluster  $x_i \in S$  belongs to.

## Third

The distance *dist* ( $x_i$ ) is the distance between  $x_i$  and its closest cluster representative object (centroid):

$$\text{dist}(x_i) = \min_{h_j \in H} L(x_i, h_j) \quad (11)$$

# Setup

## First

Set  $H$  denotes the set of cluster centroids or cluster of representative objects  $\{\mathbf{h}_1, \dots, \mathbf{h}_k\} \subset S$ .

## Second

Let  $cluster(\mathbf{x}_i)$  be the identity of the cluster  $\mathbf{x}_i \in S$  belongs to.

## Third

The distance  $dist(\mathbf{x}_i)$  is the distance between  $\mathbf{x}_i$  and its closest cluster representative object (centroid):

$$dist(\mathbf{x}_i) = \min_{\mathbf{h}_j \in H} L(\mathbf{x}_i, \mathbf{h}_j) \quad (11)$$



# Setup

## First

Set  $H$  denotes the set of cluster centroids or cluster of representative objects  $\{\mathbf{h}_1, \dots, \mathbf{h}_k\} \subset S$ .

## Second

Let  $cluster(\mathbf{x}_i)$  be the identity of the cluster  $\mathbf{x}_i \in S$  belongs to.

## Third

The distance  $dist(\mathbf{x}_i)$  is the distance between  $\mathbf{x}_i$  and its closest cluster representative object (centroid):

$$dist(\mathbf{x}_i) = \min_{\mathbf{h}_j \in H} L(\mathbf{x}_i, \mathbf{h}_j) \quad (11)$$

# The Main Idea

## Something Notable

We always assign  $x_i$  to the closest centroid. Therefore  $dist(x_i)$  is the minimum distance between  $x_i$  and any centroid.

# The Main Idea

## Something Notable

We always assign  $x_i$  to the closest centroid. Therefore  $dist(x_i)$  is the minimum distance between  $x_i$  and any centroid.

## The Algorithm is Iterative

- We generate one cluster centroid first and then add others one by one until we get  $K$  clusters.

• The set of centroids  $H$  starts with only a single centroid,  $h_1$ .

# The Main Idea

## Something Notable

We always assign  $x_i$  to the closest centroid. Therefore  $dist(x_i)$  is the minimum distance between  $x_i$  and any centroid.

## The Algorithm is Iterative

- We generate one cluster centroid first and then add others one by one until we get  $K$  clusters.
- The set of centroids  $H$  starts with only a single centroid,  $h_1$ .

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- **The  $K$ -Center Algorithm**
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# Algorithm

## Step 1

- Randomly select an object  $x_j$  from  $S$ , let  $h_1 = x_j$ ,  $H = \{h_1\}$ .
- It does not matter how  $x_j$  is selected.

# Algorithm

## Step 1

- Randomly select an object  $x_j$  from  $S$ , let  $h_1 = x_j$ ,  $H = \{h_1\}$ .
- It does not matter how  $x_j$  is selected.

## Step 2

For  $j = 1$  to  $n$ :

- $dist(x_i) = L(x_i, h_1)$ .
- $cluster(x_i)$ .

# Algorithm

## Step 1

- Randomly select an object  $x_j$  from  $S$ , let  $h_1 = x_j$ ,  $H = \{h_1\}$ .
- It does not matter how  $x_j$  is selected.

## Step 2

For  $j = 1$  to  $n$ :

- $dist(x_j) = L(x_j, h_1)$ .
- $cluster(x_j)$ .

In other words

- The  $dist(x_j)$  is the computed distance between  $x_j$  and  $h_1$ .
- Because so far we only have one cluster, we will assign a cluster label 1 to every  $x_j$ .



# Algorithm

## Step 1

- Randomly select an object  $x_j$  from  $S$ , let  $h_1 = x_j$ ,  $H = \{h_1\}$ .
- It does not matter how  $x_j$  is selected.

## Step 2

For  $j = 1$  to  $n$ :

①  $dist(x_i) = L(x_i, h_1)$ .

② *cluster* ( $x_i$ ).

in other words

- The  $dist(x_i)$  is the computed distance between  $x_j$  and  $h_1$ .
- Because so far we only have one cluster, we will assign a cluster label 1 to every  $x_j$ .

# Algorithm

## Step 1

- Randomly select an object  $x_j$  from  $S$ , let  $h_1 = x_j$ ,  $H = \{h_1\}$ .
- It does not matter how  $x_j$  is selected.

## Step 2

For  $j = 1$  to  $n$ :

- 1  $dist(x_i) = L(x_i, h_1)$ .
- 2  $cluster(x_i)$ .

in other words

- The  $dist(x_i)$  is the computed distance between  $x_j$  and  $h_1$ .
- Because so far we only have one cluster, we will assign a cluster label 1 to every  $x_j$ .

# Algorithm

## Step 1

- Randomly select an object  $x_j$  from  $S$ , let  $h_1 = x_j$ ,  $H = \{h_1\}$ .
- It does not matter how  $x_j$  is selected.

## Step 2

For  $j = 1$  to  $n$ :

- 1  $dist(x_i) = L(x_i, h_1)$ .
- 2  $cluster(x_i)$ .

## In other words

- The  $dist(x_i)$  is the computed distance between  $x_j$  and  $h_1$ .
- Because so far we only have one cluster, we will assign a cluster label 1 to every  $x_j$ .

# Algorithm

## Step 1

- Randomly select an object  $x_j$  from  $S$ , let  $h_1 = x_j$ ,  $H = \{h_1\}$ .
- It does not matter how  $x_j$  is selected.

## Step 2

For  $j = 1$  to  $n$ :

- 1  $dist(x_i) = L(x_i, h_1)$ .
- 2  $cluster(x_i)$ .

## In other words

- The  $dist(x_i)$  is the computed distance between  $x_j$  and  $h_1$ .
- Because so far we only have one cluster, we will assign a cluster label 1 to every  $x_j$ .

# Algorithm

## Step 3

For  $i = 2$  to  $K$

①  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$

② Choose  $h_i \in S - H$  such that  $\text{dist}(h_i) == D$

③  $H = H \cup \{h_i\}$

④ for  $j = 1$  to  $N$

⑤     if  $L(\mathbf{x}_j, h_i) \leq \text{dist}(\mathbf{x}_j)$

⑥          $\text{dist}(\mathbf{x}_j) = L(\mathbf{x}_j, h_i)$

⑦          $\text{cluster}(\mathbf{x}_j) = i$

# Algorithm

## Step 3

For  $i = 2$  to  $K$

- 1  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$
- 2 Choose  $\mathbf{h}_i \in S - H$  such that  $\text{dist}(\mathbf{h}_i) == D$
- 3  $H = H \cup \{\mathbf{h}_i\}$
- 4 for  $j = 1$  to  $N$
- 5     if  $L(\mathbf{x}_j, \mathbf{h}_i) \leq \text{dist}(\mathbf{x}_j)$
- 6          $\text{dist}(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
- 7          $\text{cluster}(\mathbf{x}_j) = i$

# Algorithm

## Step 3

For  $i = 2$  to  $K$

- 1  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$
  - 2 Choose  $\mathbf{h}_i \in S - H$  such that  $\text{dist}(\mathbf{h}_i) == D$
  - 3  $H = H \cup \{\mathbf{h}_i\}$
- for  $j = 1$  to  $N$
- if  $L(\mathbf{x}_j, \mathbf{h}_i) \leq \text{dist}(\mathbf{x}_j)$
  - $\text{dist}(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
  - $\text{cluster}(\mathbf{x}_j) = i$

# Algorithm

## Step 3

For  $i = 2$  to  $K$

- 1  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$
- 2 Choose  $\mathbf{h}_i \in S - H$  such that  $\text{dist}(\mathbf{h}_i) == D$
- 3  $H = H \cup \{\mathbf{h}_i\}$
- 4 for  $j = 1$  to  $N$ 
  - if  $L(\mathbf{x}_j, \mathbf{h}_i) \leq \text{dist}(\mathbf{x}_j)$
  - $\text{dist}(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
  - $\text{cluster}(\mathbf{x}_j) = i$



# Algorithm

## Step 3

For  $i = 2$  to  $K$

- 1  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$
- 2 Choose  $\mathbf{h}_i \in S - H$  such that  $\text{dist}(\mathbf{h}_i) == D$
- 3  $H = H \cup \{\mathbf{h}_i\}$
- 4 for  $j = 1$  to  $N$
- 5     if  $L(\mathbf{x}_j, \mathbf{h}_i) \leq \text{dist}(\mathbf{x}_j)$
- 6          $\text{dist}(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
- 7          $\text{cluster}(\mathbf{x}_j) = i$

# Algorithm

## Step 3

For  $i = 2$  to  $K$

- 1  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$
- 2 Choose  $\mathbf{h}_i \in S - H$  such that  $\text{dist}(\mathbf{h}_i) == D$
- 3  $H = H \cup \{\mathbf{h}_i\}$
- 4 for  $j = 1$  to  $N$
- 5     if  $L(\mathbf{x}_j, \mathbf{h}_i) \leq \text{dist}(\mathbf{x}_j)$
- 6          $\text{dist}(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
- 7      $\text{cluster}(\mathbf{x}_j) = i$

# Algorithm

## Step 3

For  $i = 2$  to  $K$

- 1  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$
- 2 Choose  $\mathbf{h}_i \in S - H$  such that  $\text{dist}(\mathbf{h}_i) == D$
- 3  $H = H \cup \{\mathbf{h}_i\}$
- 4 for  $j = 1$  to  $N$
- 5     if  $L(\mathbf{x}_j, \mathbf{h}_i) \leq \text{dist}(\mathbf{x}_j)$
- 6          $\text{dist}(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
- 7          $\text{cluster}(\mathbf{x}_j) = i$

# Thus

## As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to  $K$ .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
  - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

# Thus

## As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to  $K$ .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
  - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

## THE WORST POINT

- It is added to the set  $H$ .
- To stress gain, points already included in  $H$  are not among the consideration.

# Thus

## As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to  $K$ .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
  - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

## This worst point

- It is added to the set  $H$ .
- To stress gain, points already included in  $H$  are not among the consideration.

# Thus

## As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to  $K$ .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
  - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

## This worst point

- It is added to the set  $H$ .
- To stress gain, points already included in  $H$  are not among the consideration.

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- **Notes in Implementation**
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function



# Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
  - ▶ MakeSet
  - ▶ Find
  - ▶ Union
  - ▶ Remove iteratively through the disjoint trees.

# Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
  - ▶ MakeSet
  - ▶ Find
  - ▶ Union
  - ▶ Remove iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances

Thus, each node-element for the sets must have a field  $dist(x_j)$  such that

$$dist(x_j) = L(x_j, Find(x_j)) \quad (12)$$

# Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
  - ▶ MakeSet
  - ▶ Find
  - ▶ Union
  - ▶ Remove iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances

Thus, each node-element for the sets must have a field  $dist(x_j)$  such that

$$dist(x_j) = L(x_j, Find(x_j)) \quad (12)$$

# Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
  - ▶ MakeSet
  - ▶ Find
  - ▶ Union
  - ▶ Remove iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances

Thus, each node-element for the sets must have a field  $dist(x_j)$  such that

$$dist(x_j) = L(x_j, Find(x_j)) \quad (12)$$

# Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
  - ▶ MakeSet
  - ▶ Find
  - ▶ Union
  - ▶ Remove Iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances

Thus, each node-element for the sets must have a field  $dist(x_j)$  such that

$$dist(x_j) = L(x_j, Find(x_j)) \quad (12)$$

# Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
  - ▶ MakeSet
  - ▶ Find
  - ▶ Union
  - ▶ Remove iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances

Thus, each node-element for the sets must have a field  $dist(\mathbf{x}_j)$  such that

$$dist(\mathbf{x}_j) = L(\mathbf{x}_j, Find(\mathbf{x}_j)) \quad (12)$$

Although

Other things need to be taken in consideration

I will allow to you to think about them

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- **Examples**
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

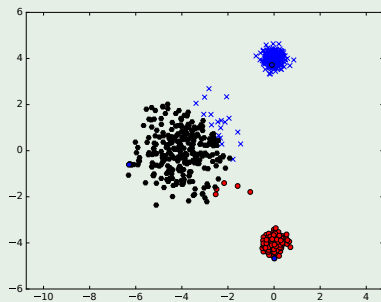
## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

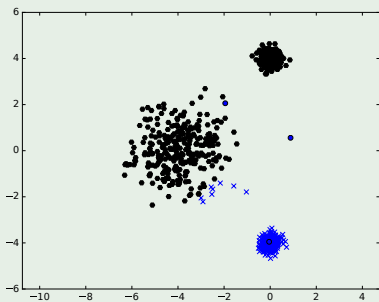


## Example

Running the k-center and k-means algorithms allows to see that for different densities k-center is more robust



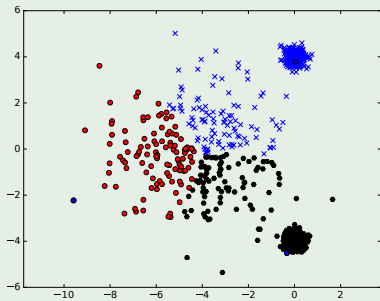
K-Center



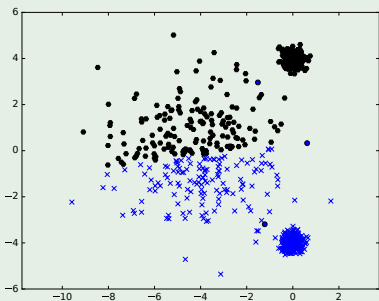
K-Means

## Example

Decreasing the density of one of the clusters, we see a degradation on the clusters



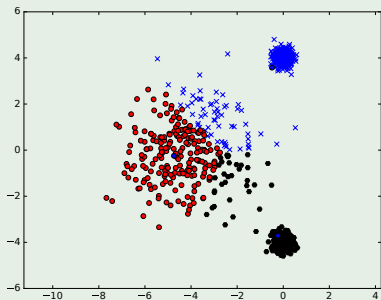
K-Center



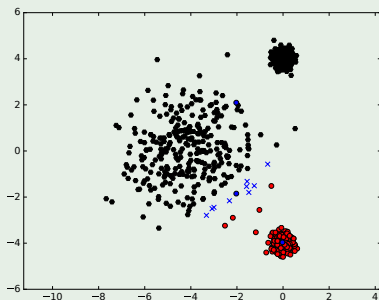
K-Means

## Using Centroids of the K-center to initialize K-mean

Thus, we can use the centroids of K-center to try to improve upon K-means to a certain degree



K-Means using the K-Center Centroids



K-Means using random centroids

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- **$K$ -Center Algorithm Properties**
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# The Running Time

## We have that

The running time of the algorithm is  $O(KN)$ , where  $K$  is the number of clusters generated and  $N$  is the size of the data set.

## Why?

Because  $K$ -center only requires pairwise distance between any point and the centroids.

# The Running Time

## We have that

The running time of the algorithm is  $O(KN)$ , where  $K$  is the number of clusters generated and  $N$  is the size of the data set.

## Why?

Because  $K$ -center only requires pairwise distance between any point and the centroids.

# Main Bound of the $K$ -center algorithm

## Lemma

Given the distance measure  $L$  satisfying the **triangle inequality**.

- If the partition obtained by the greedy algorithm is  $\tilde{S}$  and the optimal partition be  $S^*$ , such that the cluster size of  $\tilde{S}$  be  $\tilde{D}$  and the one for  $S^*$  is  $D^*$ , then

$$\tilde{D} \leq 2D^* \quad (13)$$

## Main Bound of the $K$ -center algorithm

### Lemma

Given the distance measure  $L$  satisfying the **triangle inequality**.

- If the partition obtained by the greedy algorithm is  $\tilde{S}$  and the optimal partition be  $S^*$ , such that the cluster size of  $\tilde{S}$  be  $\tilde{D}$  and the one for  $S^*$  is  $D^*$ , then

$$\tilde{D} \leq 2D^*$$

(13)



# Main Bound of the $K$ -center algorithm

## Lemma

Given the distance measure  $L$  satisfying the **triangle inequality**.

- If the partition obtained by the greedy algorithm is  $\tilde{S}$  and the optimal partition be  $S^*$ , such that the cluster size of  $\tilde{S}$  be  $\tilde{D}$  and the one for  $S^*$  is  $D^*$ , then

$$\tilde{D} \leq 2D^* \quad (13)$$

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- **$K$ -Center Algorithm proof of correctness**

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

## Proof

If we look only at the first  $j$  centroid

It generates a partition  $j$  with size  $D_j$  and also:

- The cluster size is the size of the biggest cluster in the current partition.
- The size of every cluster is defined as the maximum distance between a point and the cluster centroid

## Proof

If we look only at the first  $j$  centroid

It generates a partition  $j$  with size  $D_j$  and also:

- The cluster size is the size of the biggest cluster in the current partition.
- The size of every cluster is defined as the maximum distance between a point and the cluster centroid

Thus

$$D_1 \geq D_2 \geq D_3 \dots \quad (14)$$

# Proof

## If we look only at the first $j$ centroid

It generates a partition  $j$  with size  $D_j$  and also:

- The cluster size is the size of the biggest cluster in the current partition.
- The size of every cluster is defined as the maximum distance between a point and the cluster centroid

Thus

$$D_1 \geq D_2 \geq D_3 \dots \quad (14)$$

- Because  $D = \max_{x_j, x_j \in S-H} \text{dist}(x_j)$  and lines 5-7 in the step 3 and using induction!!!
- You can prove that part by yourselves.

## Proof

If we look only at the first  $j$  centroid

It generates a partition  $j$  with size  $D_j$  and also:

- The cluster size is the size of the biggest cluster in the current partition.
- The size of every cluster is defined as the maximum distance between a point and the cluster centroid

Thus

$$D_1 \geq D_2 \geq D_3 \dots \quad (14)$$

- Because  $D = \max_{x_j, x_j \in S-H} \text{dist}(x_j)$  and lines 5-7 in the step 3 and using induction!!!
- You can prove that part by yourselves.

## Proof

### If we look only at the first $j$ centroid

It generates a partition  $j$  with size  $D_j$  and also:

- The cluster size is the size of the biggest cluster in the current partition.
- The size of every cluster is defined as the maximum distance between a point and the cluster centroid

### Thus

$$D_1 \geq D_2 \geq D_3 \dots \quad (14)$$

### Why?

- Because  $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} \text{dist}(\mathbf{x}_j)$  and lines 5-7 in the step 3 and using induction!!!
- You can prove that part by yourselves.

## Proof

### If we look only at the first $j$ centroid

It generates a partition  $j$  with size  $D_j$  and also:

- The cluster size is the size of the biggest cluster in the current partition.
- The size of every cluster is defined as the maximum distance between a point and the cluster centroid

### Thus

$$D_1 \geq D_2 \geq D_3 \dots \quad (14)$$

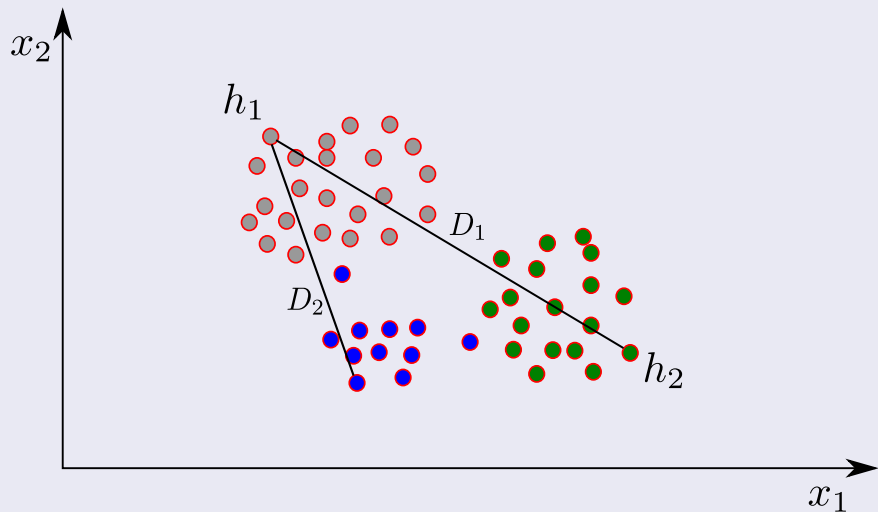
### Why?

- Because  $D = \max_{x_j: x_j \in S-H} \text{dist}(x_j)$  and lines 5-7 in the step 3 and using induction!!!
- You can prove that part by yourselves..



## Graphically

It is easy to see that when the inner loop changes distances



## Now

It is necessary to prove

$$\forall i < j, L(\mathbf{h}_i, \mathbf{h}_j) \geq D_{j-1} \quad (15)$$

Thus

$D_{j-1}$  is a lower bound for the distance between  $\mathbf{h}_i$  and  $\mathbf{h}_j$ .

## Now

It is necessary to prove

$$\forall i < j, L(\mathbf{h}_i, \mathbf{h}_j) \geq D_{j-1} \quad (15)$$

Thus

$D_{j-1}$  is a lower bound for the distance between  $\mathbf{h}_i$  and  $\mathbf{h}_j$ .

Now

It is necessary to prove

$$\forall i < j, L(\mathbf{h}_i, \mathbf{h}_j) \geq D_{j-1} \quad (15)$$

Thus

$D_{j-1}$  is a lower bound for the distance between  $\mathbf{h}_i$  and  $\mathbf{h}_j$ .

## Proof of the Previous Statement

It is possible to see that

$$L(\mathbf{h}_{j-2}, \mathbf{h}_j) \geq L(\mathbf{h}_{j-1}, \mathbf{h}_j) \quad (16)$$

Proof:

Assume it is not true. Then,  $L(\mathbf{h}_{j-2}, \mathbf{h}_j) < L(\mathbf{h}_{j-1}, \mathbf{h}_j)$

## Proof of the Previous Statement

It is possible to see that

$$L(\mathbf{h}_{j-2}, \mathbf{h}_j) \geq L(\mathbf{h}_{j-1}, \mathbf{h}_j) \quad (16)$$

How?

Assume it is not true. Then,  $L(\mathbf{h}_{j-2}, \mathbf{h}_j) < L(\mathbf{h}_{j-1}, \mathbf{h}_j)$

Then

We have that given the partition  $\tilde{S}$  generated by the greedy algorithm

- $\mathbf{h}_{j-2} \in \tilde{C}_{j-2}$  ,  $\mathbf{h}_{j-2} \in \tilde{C}_j$  for  $\tilde{C}_{j-2}, \tilde{C}_j \in \tilde{S}$
- It is a contradiction because  $\mathbf{h}_{j-2}$  is generated by the algorithm such that cannot be in any other cluster!!!

Thus iteratively

$$L(\mathbf{h}_1, \mathbf{h}_j) \geq L(\mathbf{h}_2, \mathbf{h}_j) \geq L(\mathbf{h}_3, \mathbf{h}_j) \geq \dots \geq L(\mathbf{h}_{j-1}, \mathbf{h}_j) = D_{j-1} \quad (17)$$

Then

We have that given the partition  $\tilde{\mathcal{S}}$  generated by the greedy algorithm

- $\mathbf{h}_{j-2} \in \tilde{\mathcal{C}}_{j-2}$  ,  $\mathbf{h}_{j-2} \in \tilde{\mathcal{C}}_j$  for  $\tilde{\mathcal{C}}_{j-2}, \tilde{\mathcal{C}}_j \in \tilde{\mathcal{S}}$
- It is a contradiction because  $\mathbf{h}_{j-2}$  is generated by the algorithm such that cannot be in any other cluster!!!

Thus iteratively

$$L(\mathbf{h}_1, \mathbf{h}_j) \geq L(\mathbf{h}_2, \mathbf{h}_j) \geq L(\mathbf{h}_3, \mathbf{h}_j) \geq \dots \geq L(\mathbf{h}_{j-1}, \mathbf{h}_j) = D_{j-1} \quad (17)$$



## Not only that

Not only that

$$\forall j, \exists i < j, L(h_i, h_j) = D_{j-1} \quad (18)$$

Therefore

Therefore,  $D_{j-1}$  is not only the lower bound for the distance between  $h_i$  and  $h_j$ , it is also the exact boundary for a specific  $i$ .

Not only that

Not only that

$$\forall j, \exists i < j, L(\mathbf{h}_i, \mathbf{h}_j) = D_{j-1} \quad (18)$$

Therefore

Therefore,  $D_{j-1}$  is not only the lower bound for the distance between  $\mathbf{h}_i$  and  $\mathbf{h}_j$ , it is also the exact boundary for a specific  $i$ .

Not only that

Not only that

$$\forall j, \exists i < j, L(\mathbf{h}_i, \mathbf{h}_j) = D_{j-1} \quad (18)$$

Therefore

Therefore,  $D_{j-1}$  is not only the lower bound for the distance between  $\mathbf{h}_i$  and  $\mathbf{h}_j$ , it is also the exact boundary for a specific  $i$ .

## If, we continue with the proof

### Now

- Let us consider the optimal partition  $S^*$  with  $K$  clusters and its size  $D^*$ .
  - ▶ **Suppose the greedy algorithm generates** the centroids  $\tilde{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K\}$ .
  - ▶ For the proof, we are adding one more,  $\mathbf{h}_{K+1}$ .
  - ▶ This can be done without losing generality.

## If, we continue with the proof

### Now

- Let us consider the optimal partition  $S^*$  with  $K$  clusters and its size  $D^*$ .
  - ▶ **Suppose the greedy algorithm generates** the centroids  $\tilde{H} = \{h_1, h_2, \dots, h_K\}$ .
  - ▶ For the proof, we are adding one more,  $h_{K+1}$ .
  - ▶ This can be done without losing generality.

According to the pigeonhole principle, at least two of the centroids among

$\{h_1, h_2, \dots, h_K, h_{K+1}\}$  will fall into one cluster  $k$  of the partition  $S^*$ .

If, we continue with the proof

Now

- Let us consider the optimal partition  $S^*$  with  $K$  clusters and its size  $D^*$ .
  - ▶ **Suppose the greedy algorithm generates** the centroids  $\tilde{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K\}$ .
  - ▶ For the proof, we are adding one more,  $\mathbf{h}_{K+1}$ .
  - ▶ This can be done without losing generality.

According to the pigeonhole principle, at least two of the centroids among

$\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{h}_{K+1}\}$  will fall into one cluster  $k$  of the partition  $S^*$ .

Without loss of generality,

$1 \leq i < k < j \leq K + 1 \Rightarrow$  Using the triangle inequality:

$$L(\mathbf{h}_i, \mathbf{h}_j) \leq L(\mathbf{h}_i, \mathbf{h}_k) + L(\mathbf{h}_k, \mathbf{h}_j) \leq D^* + D^* = 2D^* \quad (19)$$

If, we continue with the proof

Now

- Let us consider the optimal partition  $S^*$  with  $K$  clusters and its size  $D^*$ .
  - ▶ **Suppose the greedy algorithm generates** the centroids  $\tilde{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K\}$ .
  - ▶ For the proof, we are adding one more,  $\mathbf{h}_{K+1}$ .
  - ▶ This can be done without losing generality.

According to the pigeonhole principle, at least two of the centroids among

$\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{h}_{K+1}\}$  will fall into one cluster  $k$  of the partition  $S^*$ .

Thus assume

$1 \leq i < k < j \leq K + 1 \Rightarrow$  Using the triangle inequality:

$$L(\mathbf{h}_i, \mathbf{h}_j) \leq L(\mathbf{h}_i, \mathbf{h}_k) + L(\mathbf{h}_k, \mathbf{h}_j) \leq D^* + D^* = 2D^* \quad (19)$$

## If, we continue with the proof

### Now

- Let us consider the optimal partition  $S^*$  with  $K$  clusters and its size  $D^*$ .
  - Suppose the greedy algorithm generates the centroids  $\tilde{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K\}$ .
  - For the proof, we are adding one more,  $\mathbf{h}_{K+1}$ .
  - This can be done without losing generality.

According to the pigeonhole principle, at least two of the centroids among

$\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{h}_{K+1}\}$  will fall into one cluster  $k$  of the partition  $S^*$ .

Thus assume

$1 \leq i < k < j \leq K + 1 \Rightarrow$  Using the triangle inequality:

$$L(\mathbf{h}_i, \mathbf{h}_j) \leq L(\mathbf{h}_i, \mathbf{h}_k) + L(\mathbf{h}_k, \mathbf{h}_j) \leq D^* + D^* = 2D^* \quad (19)$$



Then

In addition

Also  $L(\mathbf{h}_i, \mathbf{h}_j) \geq D_{j-1} \geq D_k$  then  $D_k \leq 2D^*$

Given  $S$ , the partition generated by the greedy algorithm

We define  $\Delta$  as

$$\Delta = \max_{x_j: x_j \in S} \min_{h_k: h_k \in \Pi} L(x_j, h_k) \quad (20)$$

Essentially

The maximum of all points that are not centroids that minimize the distance to some centroid for the partition generated by the greedy algorithm.

Then

In addition

Also  $L(\mathbf{h}_i, \mathbf{h}_j) \geq D_{j-1} \geq D_k$  then  $D_k \leq 2D^*$

Given  $\tilde{S}$ , the partition generated by the greedy algorithm

We define  $\Delta$  as

$$\Delta = \max_{\mathbf{x}_j: \mathbf{x}_j \in \tilde{S}} \min_{\tilde{H} \mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{x}_j, \mathbf{h}_k) \quad (20)$$

Finally

The maximum of all points that are not centroids that minimize the distance to some centroid for the partition generated by the greedy algorithm.

# Then

In addition

Also  $L(\mathbf{h}_i, \mathbf{h}_j) \geq D_{j-1} \geq D_k$  then  $D_k \leq 2D^*$

Given  $\tilde{S}$ , the partition generated by the greedy algorithm

We define  $\Delta$  as

$$\Delta = \max_{\mathbf{x}_j: \mathbf{x}_j \in \tilde{S}} \min_{\tilde{H} \mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{x}_j, \mathbf{h}_k) \quad (20)$$

Basically

The maximum of all points that are not centroids that minimize the distance to some centroid for the partition generated by the greedy algorithm.

Now, we are ready for the final part

Let  $\mathbf{h}_{K+1}$  be an element in  $\tilde{S} - \tilde{H}$

Such that

$$\min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) = \Delta \quad (21)$$

By definition

$$L(\mathbf{h}_{K+1}, \mathbf{h}_k) \geq \Delta, \quad \forall k = 1, \dots, K \quad (22)$$

Thus, we have the following sets

Let  $H_k = \{\mathbf{h}_1, \dots, \mathbf{h}_k\}$  with  $k = 1, 2, \dots, K$ .

Now, we are ready for the final part

Let  $\mathbf{h}_{K+1}$  be an element in  $\tilde{S} - \tilde{H}$

Such that

$$\min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) = \Delta \quad (21)$$

By definition

$$L(\mathbf{h}_{K+1}, \mathbf{h}_k) \geq \Delta, \quad \forall k = 1, \dots, K \quad (22)$$

Thus, we have the following sets

Let  $H_k = \{\mathbf{h}_1, \dots, \mathbf{h}_k\}$  with  $k = 1, 2, \dots, K$ .

Now, we are ready for the final part

Let  $\mathbf{h}_{K+1}$  be an element in  $\tilde{S} - \tilde{H}$

Such that

$$\min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) = \Delta \quad (21)$$

By definition

$$L(\mathbf{h}_{K+1}, \mathbf{h}_k) \geq \Delta, \quad \forall k = 1, \dots, K \quad (22)$$

Thus, we have the following sets

Let  $H_k = \{\mathbf{h}_1, \dots, \mathbf{h}_k\}$  with  $k = 1, 2, \dots, K$ .

Consider the distance between  $\mathbf{h}_i$  and  $\mathbf{h}_j$  for  $i < j \leq K$

- According the greedy algorithm

$$\min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_j, \mathbf{h}_k) \geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{x}_l, \mathbf{h}_k) \text{ for any } \mathbf{x}_l \in \tilde{S} - H_j$$

- Basically remember that the  $\mathbf{h}_i$  are obtained by finding the farthest points.

Now, we have

Since  $\mathbf{h}_{K+1} \in \tilde{S} - \tilde{H}$  and  $\tilde{S} - \tilde{H} \subset \tilde{S} - H_j$

$$\begin{aligned} L(\mathbf{h}_j, \mathbf{h}_i) &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_j, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &= \Delta \end{aligned}$$



Now, we have

Since  $\mathbf{h}_{K+1} \in \tilde{S} - \tilde{H}$  and  $\tilde{S} - \tilde{H} \subset \tilde{S} - H_j$

$$\begin{aligned} L(\mathbf{h}_j, \mathbf{h}_i) &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_j, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &= \Delta \end{aligned}$$

We have shown that for any for  $j = 1, \dots, K+1$

$$L(\mathbf{h}_j, \mathbf{h}_i) \geq \Delta \quad (23)$$

Now, we have

Since  $\mathbf{h}_{K+1} \in \tilde{S} - \tilde{H}$  and  $\tilde{S} - \tilde{H} \subset \tilde{S} - H_j$

$$\begin{aligned} L(\mathbf{h}_j, \mathbf{h}_i) &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_j, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &=\Delta \end{aligned}$$

We have shown that for any  $i, j \leq K+1$

$$L(\mathbf{h}_j, \mathbf{h}_i) \geq \Delta \quad (23)$$

Now, we have

Since  $\mathbf{h}_{K+1} \in \tilde{S} - \tilde{H}$  and  $\tilde{S} - \tilde{H} \subset \tilde{S} - H_j$

$$\begin{aligned} L(\mathbf{h}_j, \mathbf{h}_i) &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_j, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &= \Delta \end{aligned}$$

We have shown that for any for  $i < j \leq K + 1$

$$L(\mathbf{h}_j, \mathbf{h}_i) \geq \Delta \quad (23)$$

Now, we have

Since  $\mathbf{h}_{K+1} \in \tilde{S} - \tilde{H}$  and  $\tilde{S} - \tilde{H} \subset \tilde{S} - H_j$

$$\begin{aligned} L(\mathbf{h}_j, \mathbf{h}_i) &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_j, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in H_{j-1}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &\geq \min_{\mathbf{h}_k: \mathbf{h}_k \in \tilde{H}} L(\mathbf{h}_{K+1}, \mathbf{h}_k) \\ &= \Delta \end{aligned}$$

We have shown that for any for  $i < j \leq K + 1$

$$L(\mathbf{h}_j, \mathbf{h}_i) \geq \Delta \tag{23}$$

## Now

Consider the optimal partition  $S^* = \{C_1^*, C_2^*, \dots, C_K^*\}$

Thus at least 2 of the  $K + 1$  elements  $h_1, h_2, \dots, h_{K+1}$  will be covered by one cluster.

Assume that

$h_i$  and  $h_j$  belong to the same cluster in  $S^*$ . Then  $L(h_i, h_j) \leq D^*$ .

In addition

We have that since  $L(h_i, h_j) \geq \Delta$  then  $\Delta \leq D^*$

## Now

Consider the optimal partition  $S^* = \{C_1^*, C_2^*, \dots, C_K^*\}$

Thus at least 2 of the  $K + 1$  elements  $h_1, h_2, \dots, h_{K+1}$  will be covered by one cluster.

Assume that

$h_i$  and  $h_j$  belong to the same cluster in  $S^*$ . Then  $L(h_i, h_j) \leq D^*$ .

In addition

We have that since  $L(h_i, h_j) \geq \Delta$  then  $\Delta \leq D^*$

## Now

Consider the optimal partition  $S^* = \{C_1^*, C_2^*, \dots, C_K^*\}$

Thus at least 2 of the  $K + 1$  elements  $h_1, h_2, \dots, h_{K+1}$  will be covered by one cluster.

Assume that

$h_i$  and  $h_j$  belong to the same cluster in  $S^*$ . Then  $L(h_i, h_j) \leq D^*$ .

In addition

We have that since  $L(h_i, h_j) \geq \Delta$  then  $\Delta \leq D^*$

## In addition

Consider elements  $\mathbf{x}_m$  and  $\mathbf{x}_n$  in any cluster represented by  $\mathbf{h}_k$

$$L(\mathbf{x}_m, \mathbf{h}_k) \leq \Delta \text{ and } L(\mathbf{x}_n, \mathbf{h}_k) \leq \Delta \quad (24)$$

By Triangle Inequality

$$L(\mathbf{x}_m, \mathbf{x}_n) \leq L(\mathbf{x}_m, \mathbf{h}_k) + L(\mathbf{x}_n, \mathbf{h}_k) \leq 2\Delta \quad (25)$$

Finally, there are two elements  $\mathbf{x}_m$  and  $\mathbf{x}_n$  in a cluster such that  $D = L(\mathbf{x}_m, \mathbf{x}_n)$

$$\tilde{D} = \max_k D_k \leq 2\Delta \leq 2D^* \quad (26)$$



## In addition

Consider elements  $\mathbf{x}_m$  and  $\mathbf{x}_n$  in any cluster represented by  $\mathbf{h}_k$

$$L(\mathbf{x}_m, \mathbf{h}_k) \leq \Delta \text{ and } L(\mathbf{x}_n, \mathbf{h}_k) \leq \Delta \quad (24)$$

By Triangle Inequality

$$L(\mathbf{x}_m, \mathbf{x}_k) \leq L(\mathbf{x}_m, \mathbf{h}_k) + L(\mathbf{x}_n, \mathbf{h}_k) \leq 2\Delta \quad (25)$$

Finally, there are two elements  $\mathbf{x}_m$  and  $\mathbf{x}_n$  in a cluster such that  $D = L(\mathbf{x}_m, \mathbf{x}_n)$

$$\tilde{D} = \max_k D_k \leq 2\Delta \leq 2D^* \quad (26)$$

In addition

Consider elements  $\mathbf{x}_m$  and  $\mathbf{x}_n$  in any cluster represented by  $\mathbf{h}_k$

$$L(\mathbf{x}_m, \mathbf{h}_k) \leq \Delta \text{ and } L(\mathbf{x}_n, \mathbf{h}_k) \leq \Delta \quad (24)$$

By Triangle Inequality

$$L(\mathbf{x}_m, \mathbf{x}_n) \leq L(\mathbf{x}_m, \mathbf{h}_k) + L(\mathbf{x}_n, \mathbf{h}_k) \leq 2\Delta \quad (25)$$

Finally, there are two elements  $\mathbf{x}_m$  and  $\mathbf{x}_n$  in a cluster such that  $\tilde{D} = L(\mathbf{x}_m, \mathbf{x}_n)$

$$\tilde{D} = \max_k D_k \leq 2\Delta \leq 2D^* \quad (26)$$

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- **Fuzzy Clustering**
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# Some of the Fuzzy Clustering Models

## Fuzzy Clustering Model

Bezdek, 1981

## Possibilistic Clustering Model

Krishnapuram - Keller, 1993

## Fuzzy Possibilistic Clustering Model

N. Pal - K. Pal - Bezdek, 1997

# Some of the Fuzzy Clustering Models

## Fuzzy Clustering Model

Bezdek, 1981

## Possibilistic Clustering Model

Krishnapuram - Keller, 1993

## Fuzzy Possibilistic Clustering Model

N. Pal - K. Pal - Bezdek, 1997

# Some of the Fuzzy Clustering Models

## Fuzzy Clustering Model

Bezdek, 1981

## Possibilistic Clustering Model

Krishnapuram - Keller, 1993

## Fuzzy Possibilistic Clustering Model

N. Pal - K. Pal - Bezdek, 1997

# Fuzzy $C$ -Means Clustering

## The input an unlabeled data set

- $X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N\}$ .
- $\mathbf{x}_k \in \mathbb{R}^p$

## Output

- A partition  $\mathcal{S}$  of the  $X$  as a matrix  $U$  of  $C \times N$ .
- Set of cluster centers  $V = \{v_1, v_2, \dots, v_C\} \subset \mathbb{R}^p$

# Fuzzy $C$ -Means Clustering

## The input an unlabeled data set

- $X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N\}$ .
- $\mathbf{x}_k \in \mathbb{R}^p$

## Output

- A partition  $\mathcal{S}$  of the  $X$  as a matrix  $U$  of  $C \times N$ .
- Set of cluster centers  $V = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_C\} \subset \mathbb{R}^p$



# What we want

## Creation of the Cost Function

First:

- We can use a distance defined as:

$$\|\mathbf{x}_k - \mathbf{v}_i\| = \sqrt{(\mathbf{x}_k - \mathbf{v}_i)^T (\mathbf{x}_k - \mathbf{v}_i)} \quad (27)$$

The euclidean distance from a point  $k$  to a centroid  $i$ .

NOTE other distances based in Mahalonobis can be taken in consideration.

# What we want

## Creation of the Cost Function

First:

- We can use a distance defined as:

$$\|\mathbf{x}_k - \mathbf{v}_i\| = \sqrt{(\mathbf{x}_k - \mathbf{v}_i)^T (\mathbf{x}_k - \mathbf{v}_i)} \quad (27)$$

The euclidean distance from a point  $k$  to a centroid  $i$ .

NOTE other distances based in Mahalonobis can be taken in consideration.

# What we want

## Creation of the Cost Function

First:

- We can use a distance defined as:

$$\|\mathbf{x}_k - \mathbf{v}_i\| = \sqrt{(\mathbf{x}_k - \mathbf{v}_i)^T (\mathbf{x}_k - \mathbf{v}_i)} \quad (27)$$

The euclidean distance from a point  $k$  to a centroid  $i$ .

**NOTE other distances based in Mahalonobis can be taken in consideration.**

Do you remember the cost function for  $K$ -means?

Finding a partition  $S$  that minimizes the following function

$$\min_S \sum_{k=1}^N \sum_{k:\mathbf{x}_k \in C_i} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (28)$$

Where  $\mathbf{v}_i = \frac{1}{N_i} \sum_{\mathbf{x}_k \in C_i} \mathbf{x}_k$

We can rewrite the previous equation as

$$\min_S \sum_{k=1}^N \sum_{i=1}^C I(\mathbf{x}_k \in C_i) \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (29)$$

Do you remember the cost function for  $K$ -means?

Finding a partition  $S$  that minimizes the following function

$$\min_S \sum_{k=1}^N \sum_{k:\mathbf{x}_k \in C_i} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (28)$$

Where  $\mathbf{v}_i = \frac{1}{N_i} \sum_{\mathbf{x}_k \in C_i} \mathbf{x}_k$

We can rewrite the previous equation as

$$\min_S \sum_{k=1}^N \sum_{i=1}^C I(\mathbf{x}_k \in C_i) \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (29)$$

In addition

Did you notice that the membership is always one or zero?

$$\min_S \sum_{k=1}^N \sum_{i=1}^C \overbrace{I(x_k \in C_i)}^{\text{Membership}} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (30)$$

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
    - Using the Lagrange Multipliers
    - Examples
    - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

Thus, we can rethink the membership using something  
“Fuzzy”

What if we modify the cost function to something like this

$$\min_S \sum_{k=1}^N \sum_{i=1}^C \overbrace{\text{Membership}}^{\text{Fuzzy Value}} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (31)$$

This means that we think that each cluster  $C_i$  is “Fuzzy”

We can assume a fuzzy set for the cluster  $C_i$  with membership function:



Thus, we can rethink the membership using something “Fuzzy”

What if we modify the cost function to something like this

$$\min_S \sum_{k=1}^N \sum_{i=1}^C \overbrace{\text{Membership}}^{\text{Fuzzy Value}} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (31)$$

This means that we think that each cluster  $C_i$  is “Fuzzy”

We can assume a fuzzy set for the cluster  $C_i$  with membership function:

$$A_i : \mathbb{R}^p \rightarrow [0, 1] \quad (32)$$

Such that we can tune it by using a power i.e. decreasing it by a  $m$  power.

Thus, we can rethink the membership using something “Fuzzy”

What if we modify the cost function to something like this

$$\min_S \sum_{k=1}^N \sum_{i=1}^C \overbrace{\text{Membership Fuzzy Value}} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (31)$$

This means that we think that each cluster  $C_i$  is “Fuzzy”

We can assume a fuzzy set for the cluster  $C_i$  with membership function:

$$A_i : \mathbb{R}^p \rightarrow [0, 1] \quad (32)$$

Such that we can tune it by using a power i.e. decreasing it by a  $m$  power.

Thus, we can rethink the membership using something “Fuzzy”

What if we modify the cost function to something like this

$$\min_S \sum_{k=1}^N \sum_{i=1}^C \overbrace{\text{Membership Fuzzy Value}} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (31)$$

This means that we think that each cluster  $C_i$  is “Fuzzy”

We can assume a fuzzy set for the cluster  $C_i$  with membership function:

$$A_i : \mathbb{R}^p \rightarrow [0, 1] \quad (32)$$

Such that we can tune it by using a power i.e. decreasing it by a  $m$  power.

## Under the following constraints

First

$$A_i(\mathbf{x}_k) \in [0, 1] \quad \forall i, k \quad (33)$$

Second

$$0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < N \quad \forall i \quad (34)$$

Third

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = 1 \quad \forall k \quad (35)$$

## Under the following constraints

First

$$A_i(\mathbf{x}_k) \in [0, 1] \quad \forall i, k \quad (33)$$

Second

$$0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < N \quad \forall i \quad (34)$$

Third

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = 1 \quad \forall k \quad (35)$$

## Under the following constraints

First

$$A_i(\mathbf{x}_k) \in [0, 1] \quad \forall i, k \quad (33)$$

Second

$$0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < N \quad \forall i \quad (34)$$

Third

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = 1 \quad \forall k \quad (35)$$

# Final Cost Function

## Properties

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (36)$$

# Final Cost Function

## Properties

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (36)$$

## Under the constraints

- $A_i(\mathbf{x}_k) \in [0, 1]$ , for  $1 \leq k \leq N$  and  $1 \leq i \leq C$ .
- $\sum_{i=1}^C A_i(\mathbf{x}_k) = 1$ , for  $1 \leq k \leq N$ .
- $0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < n$ , for  $1 \leq i \leq C$ .
- $m > 1$ .



# Final Cost Function

## Properties

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (36)$$

## Under the constraints

- $A_i(\mathbf{x}_k) \in [0, 1]$ , for  $1 \leq k \leq N$  and  $1 \leq i \leq C$ .
- $\sum_{i=1}^C A_i(\mathbf{x}_k) = 1$ , for  $1 \leq k \leq N$ .

- $0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < n$ , for  $1 \leq i \leq C$ .

- $m > 1$ .

# Final Cost Function

## Properties

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (36)$$

## Under the constraints

- $A_i(\mathbf{x}_k) \in [0, 1]$ , for  $1 \leq k \leq N$  and  $1 \leq i \leq C$ .
- $\sum_{i=1}^C A_i(\mathbf{x}_k) = 1$ , for  $1 \leq k \leq N$ .
- $0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < n$ , for  $1 \leq i \leq C$ .

•  $m > 1$ .

# Final Cost Function

## Properties

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (36)$$

## Under the constraints

- $A_i(\mathbf{x}_k) \in [0, 1]$ , for  $1 \leq k \leq N$  and  $1 \leq i \leq C$ .
- $\sum_{i=1}^C A_i(\mathbf{x}_k) = 1$ , for  $1 \leq k \leq N$ .
- $0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < n$ , for  $1 \leq i \leq C$ .
- $m > 1$ .

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - **Using the Lagrange Multipliers**
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

# Using the Lagrange Multipliers

## New cost function

$$\bar{J}_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 - \sum_{k=1}^N \lambda_k \left[ \sum_{i=1}^C A_i(\mathbf{x}_k) - 1 \right] \quad (37)$$

Derive with respect to  $A_i(\mathbf{x}_k)$

$$\frac{\partial \bar{J}_m(\mathcal{S})}{\partial A_i(\mathbf{x}_k)} = m A_i(\mathbf{x}_k)^{m-1} \|\mathbf{x}_k - \mathbf{v}_i\|^2 - \lambda_k = 0 \quad (38)$$

This

$$A_i(\mathbf{x}_k) = \left[ \frac{\lambda_k}{m \|\mathbf{x}_k - \mathbf{v}_i\|^2} \right]^{\frac{1}{m-1}} \quad (39)$$

## Using the Lagrange Multipliers

### New cost function

$$\bar{J}_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 - \sum_{k=1}^N \lambda_k \left[ \sum_{i=1}^C A_i(\mathbf{x}_k) - 1 \right] \quad (37)$$

### Derive with respect to $A_i(\mathbf{x}_k)$

$$\frac{\partial \bar{J}_m(\mathcal{S})}{\partial A_i(\mathbf{x}_k)} = m A_i(\mathbf{x}_k)^{m-1} \|\mathbf{x}_k - \mathbf{v}_i\|^2 - \lambda_k = 0 \quad (38)$$

This

$$A_i(\mathbf{x}_k) = \left[ \frac{\lambda_k}{m \|\mathbf{x}_k - \mathbf{v}_i\|^2} \right]^{\frac{1}{m-1}} \quad (39)$$

## Using the Lagrange Multipliers

### New cost function

$$\bar{J}_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 - \sum_{k=1}^N \lambda_k \left[ \sum_{i=1}^C A_i(\mathbf{x}_k) - 1 \right] \quad (37)$$

### Derive with respect to $A_i(\mathbf{x}_k)$

$$\frac{\partial \bar{J}_m(\mathcal{S})}{\partial A_i(\mathbf{x}_k)} = mA_i(\mathbf{x}_k)^{m-1} \|\mathbf{x}_k - \mathbf{v}_i\|^2 - \lambda_k = 0 \quad (38)$$

Thus

$$A_i(\mathbf{x}_k) = \left[ \frac{\lambda_k}{m \|\mathbf{x}_k - \mathbf{v}_i\|^2} \right]^{\frac{1}{m-1}} \quad (39)$$

## Using the Lagrange Multipliers

Sum over all  $i$ 's

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = \frac{\lambda_k^{\frac{1}{m-1}}}{m^{\frac{1}{m-1}} \|\mathbf{x}_k - \mathbf{v}_i\|^{\frac{2}{m-1}}} \quad (40)$$

Thus

$$\lambda_k = \frac{m}{\left[ \sum_{i=1}^C \frac{1}{\|\mathbf{x}_k - \mathbf{v}_i\|^{\frac{2}{m-1}}} \right]^{m-1}} \quad (41)$$

Plug Back on equation 38 using  $\lambda$  instead of  $\lambda_k$

$$\frac{m}{\left[ \sum_{j=1}^C \frac{1}{\|\mathbf{x}_k - \mathbf{v}_j\|^{\frac{2}{m-1}}} \right]^{m-1}} = m A_i(\mathbf{x}_k)^{m-1} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (42)$$



## Using the Lagrange Multipliers

Sum over all  $i$ 's

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = \frac{\lambda_k^{\frac{1}{m-1}}}{m^{\frac{1}{m-1}} \|\mathbf{x}_k - \mathbf{v}_i\|^{\frac{2}{m-1}}} \quad (40)$$

Thus

$$\lambda_k = \frac{m}{\left[ \sum_{i=1}^C \frac{1}{\|\mathbf{x}_k - \mathbf{v}_i\|^{\frac{2}{m-1}}} \right]^{m-1}} \quad (41)$$

Plug back on equation 40 using  $\lambda_k$  instead of  $\lambda_k^{\frac{1}{m-1}}$

$$\left[ \sum_{j=1}^C \frac{1}{\|\mathbf{x}_k - \mathbf{v}_j\|^{\frac{2}{m-1}}} \right]^{m-1} = m A_i(\mathbf{x}_k)^{m-1} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (42)$$

## Using the Lagrange Multipliers

Sum over all  $i$ 's

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = \frac{\lambda_k^{\frac{1}{m-1}}}{m^{\frac{1}{m-1}} \|\mathbf{x}_k - \mathbf{v}_i\|^{\frac{2}{m-1}}} \quad (40)$$

Thus

$$\lambda_k = \frac{m}{\left[ \sum_{i=1}^C \frac{1}{\|\mathbf{x}_k - \mathbf{v}_i\|^{\frac{2}{m-1}}} \right]^{m-1}} \quad (41)$$

Plug Back on equation 38 using  $j$  instead of  $i$

$$\frac{m}{\left[ \sum_{j=1}^C \frac{1}{\|\mathbf{x}_k - \mathbf{v}_j\|^{\frac{2}{m-1}}} \right]^{m-1}} = m A_i(\mathbf{x}_k)^{m-1} \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (42)$$

Finally

We have that

$$A_i(\mathbf{x}_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|\mathbf{x}_k - \mathbf{v}_i\|^2}{\|\mathbf{x}_k - \mathbf{v}_j\|^2} \right\}^{\frac{1}{m-1}} \right]} \quad (43)$$

In a similar way we have

$$\mathbf{v}_i = \frac{\sum_{k=1}^N A_i(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i(\mathbf{x}_k)^m} \quad (44)$$

Finally

We have that

$$A_i(\mathbf{x}_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|\mathbf{x}_k - \mathbf{v}_i\|^2}{\|\mathbf{x}_k - \mathbf{v}_j\|^2} \right\}^{\frac{1}{m-1}} \right]} \quad (43)$$

In a similar way we have

$$\mathbf{v}_i = \frac{\sum_{k=1}^N A_i(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i(\mathbf{x}_k)^m} \quad (44)$$

# Final Algorithm

## Fuzzy c-means

- 1 Let  $t = 0$ . Select an initial fuzzy pseudo-partition.
- 2 Calculate the initial C cluster centers using,  $v_i^{(t)} = \frac{\sum_{k=1}^N A_i^{(t)}(x_k)^m x_k}{\sum_{k=1}^N A_i^{(t)}(x_k)^m}$ .
- 3 Update for each  $x_k$  the membership function by
  - ▶ Case I:  $\|x_k - v_i^{(t)}\|^2 > 0$  for all  $i \in \{1, 2, \dots, C\}$  then
$$A_i^{(t+1)}(x_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|x_k - v_i^{(t)}\|^2}{\|x_k - v_j^{(t)}\|^2} \right\}^{\frac{1}{m-1}} \right]}$$
  - ▶ Case II:  $\|x_k - v_i^{(t)}\|^2 = 0$  for some  $i \in I \subseteq \{1, 2, \dots, C\}$  then define  $A_i^{(t+1)}(x_k)$  by any nonnegative number such that  $\sum_{i \in I} A_i^{(t+1)}(x_k) = 1$  and  $A_i^{(t+1)}(x_k) = 0$  for  $i \notin I$ .
- 4 If  $|\mathcal{S}^{(t+1)} - \mathcal{S}^{(t)}| = \max_{i,k} |A_i^{(t+1)}(x_k) - A_i^{(t)}(x_k)| \leq \epsilon$  stop; otherwise increase  $t$  and go to step 2.

# Final Algorithm

## Fuzzy c-means

1 Let  $t = 0$ . Select an initial fuzzy pseudo-partition.

2 Calculate the initial C cluster centers using,  $v_i^{(t)} = \frac{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m}$ .

3 Update for each  $x_k$  the membership function by

▶ Case I:  $\|x_k - v_i^{(t)}\|^2 > 0$  for all  $i \in \{1, 2, \dots, C\}$  then

$$A_i^{(t+1)}(x_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|x_k - v_i^{(t)}\|^2}{\|x_k - v_j^{(t)}\|^2} \right\}^{\frac{1}{m-1}} \right]}$$

▶ Case II:  $\|x_k - v_i^{(t)}\|^2 = 0$  for some  $i \in I \subseteq \{1, 2, \dots, C\}$  then define

$$A_i^{(t+1)}(x_k) \text{ by any nonnegative number such that } \sum_{i \in I} A_i^{(t+1)}(x_k) = 1 \text{ and } A_i^{(t+1)}(x_k) = 0 \text{ for } i \notin I.$$

4 If  $|\mathcal{S}^{(t+1)} - \mathcal{S}^{(t)}| = \max_{i,k} |A_i^{(t+1)}(x_k) - A_i^{(t)}(x_k)| \leq \epsilon$  stop; otherwise increase  $t$  and go to step 2.

# Final Algorithm

## Fuzzy c-means

① Let  $t = 0$ . Select an initial fuzzy pseudo-partition.

② Calculate the initial C cluster centers using,  $v_i^{(t)} = \frac{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m}$ .

③ Update for each  $x_k$  the membership function by

▶ Case I:  $\|\mathbf{x}_k - v_i^{(t)}\|^2 > 0$  for all  $i \in \{1, 2, \dots, C\}$  then

$$A_i^{(t+1)}(\mathbf{x}_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|\mathbf{x}_k - v_j^{(t)}\|^2}{\|\mathbf{x}_k - v_i^{(t)}\|^2} \right\}^{\frac{1}{m-1}} \right]}$$

▶ Case II:  $\|\mathbf{x}_k - v_i^{(t)}\|^2 = 0$  for some  $i \in I \subseteq \{1, 2, \dots, C\}$  then define

$A_i^{(t+1)}(\mathbf{x}_k)$  by any nonnegative number such that

$\sum_{i \in I} A_i^{(t+1)}(\mathbf{x}_k) = 1$  and  $A_i^{(t+1)}(\mathbf{x}_k) = 0$  for  $i \notin I$ .

④ If  $|\mathcal{S}^{(t+1)} - \mathcal{S}^{(t)}| = \max_{i,k} |A_i^{(t+1)}(\mathbf{x}_k) - A_i^{(t)}(\mathbf{x}_k)| \leq \epsilon$  stop; otherwise increase  $t$  and go to step 2.

# Final Algorithm

## Fuzzy c-means

① Let  $t = 0$ . Select an initial fuzzy pseudo-partition.

② Calculate the initial C cluster centers using,  $v_i^{(t)} = \frac{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m}$ .

③ Update for each  $x_k$  the membership function by

▶ Case I:  $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 > 0$  for all  $i \in \{1, 2, \dots, C\}$  then

$$A_i^{(t+1)}(\mathbf{x}_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2}{\|\mathbf{x}_k - \mathbf{v}_j^{(t)}\|^2} \right\}^{\frac{1}{m-1}} \right]}$$

▶ Case II:  $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 = 0$  for some  $i \in I \subseteq \{1, 2, \dots, C\}$  then define

$A_i^{(t+1)}(\mathbf{x}_k)$  by any nonnegative number such that

$\sum_{i \in I} A_i^{(t+1)}(\mathbf{x}_k) = 1$  and  $A_i^{(t+1)}(\mathbf{x}_k) = 0$  for  $i \notin I$ .

④ If  $|\mathcal{S}^{(t+1)} - \mathcal{S}^{(t)}| = \max_{i,k} |A_i^{(t+1)}(\mathbf{x}_k) - A_i^{(t)}(\mathbf{x}_k)| \leq \epsilon$  stop; otherwise increase  $t$  and go to step 2.



# Final Algorithm

## Fuzzy c-means

① Let  $t = 0$ . Select an initial fuzzy pseudo-partition.

② Calculate the initial C cluster centers using,  $v_i^{(t)} = \frac{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m}$ .

③ Update for each  $x_k$  the membership function by

▶ Case I:  $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 > 0$  for all  $i \in \{1, 2, \dots, C\}$  then

$$A_i^{(t+1)}(\mathbf{x}_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2}{\|\mathbf{x}_k - \mathbf{v}_j^{(t)}\|^2} \right\}^{\frac{1}{m-1}} \right]}$$

▶ Case II:  $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 = 0$  for some  $i \in I \subseteq \{1, 2, \dots, C\}$  then define

$A_i^{(t+1)}(\mathbf{x}_k)$  by any nonnegative number such that

$$\sum_{i \in I} A_i^{(t+1)}(\mathbf{x}_k) = 1 \text{ and } A_i^{(t+1)}(\mathbf{x}_k) = 0 \text{ for } i \notin I.$$

④ If  $|S^{(t+1)} - S^{(t)}| = \max_k |A_i^{(t+1)}(\mathbf{x}_k) - A_i^{(t)}(\mathbf{x}_k)| \leq \epsilon$  stop; otherwise increase  $t$  and go to step 2.

# Final Algorithm

## Fuzzy c-means

- 1 Let  $t = 0$ . Select an initial fuzzy pseudo-partition.
- 2 Calculate the initial C cluster centers using,  $v_i^{(t)} = \frac{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m}$ .
- 3 Update for each  $x_k$  the membership function by
  - ▶ Case I:  $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 > 0$  for all  $i \in \{1, 2, \dots, C\}$  then
$$A_i^{(t+1)}(\mathbf{x}_k) = \frac{1}{\left[ \sum_{j=1}^C \left\{ \frac{\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2}{\|\mathbf{x}_k - \mathbf{v}_j^{(t)}\|^2} \right\}^{\frac{1}{m-1}} \right]}$$
  - ▶ Case II:  $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 = 0$  for some  $i \in I \subseteq \{1, 2, \dots, C\}$  then define  $A_i^{(t+1)}(\mathbf{x}_k)$  by any nonnegative number such that  $\sum_{i \in I} A_i^{(t+1)}(\mathbf{x}_k) = 1$  and  $A_i^{(t+1)}(\mathbf{x}_k) = 0$  for  $i \notin I$ .
- 4 If  $|\mathcal{S}^{(t+1)} - \mathcal{S}^{(t)}| = \max_{i,k} |A_i^{(t+1)}(\mathbf{x}_k) - A_i^{(t)}(\mathbf{x}_k)| \leq \epsilon$  stop; otherwise increase  $t$  and go to step 2.

# Final Output

## The Matrix $U$

The elements of  $U$  are  $U_{ik} = A_i(\mathbf{x}_k)$ .

The centroids

$$V = \{v_1, v_2, \dots, v_C\}$$

# Final Output

## The Matrix $U$

The elements of  $U$  are  $U_{ik} = A_i(\mathbf{x}_k)$ .

## The centroids

$$V = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_C\}$$

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

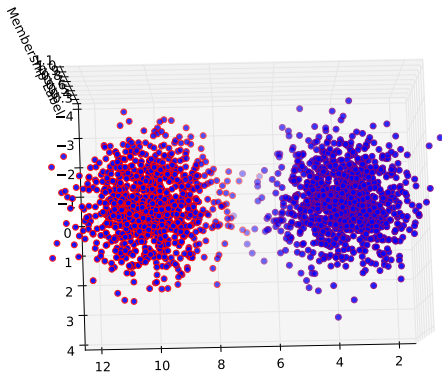
- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- **Fuzzy Clustering**
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - **Examples**
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

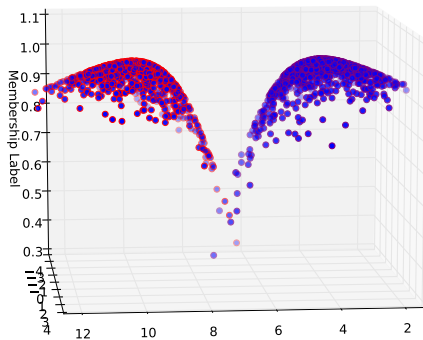
## Example

Here the clustering of two Gaussian Clusters with  $\mu_1 = (4, 0)^T$ ,  $\mu_2 = (10, 0)^T$  and variance 1.0



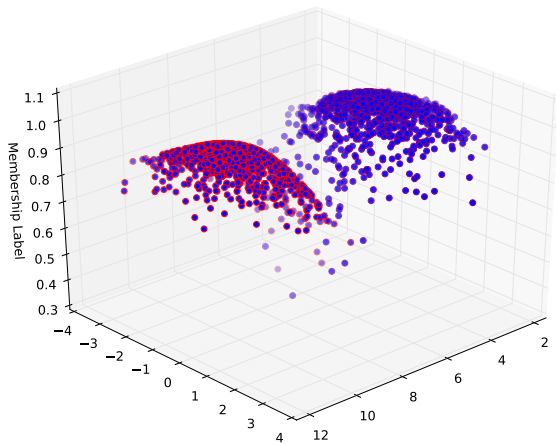
# Example

Here the clustering of two Gaussian Clusters



# Example

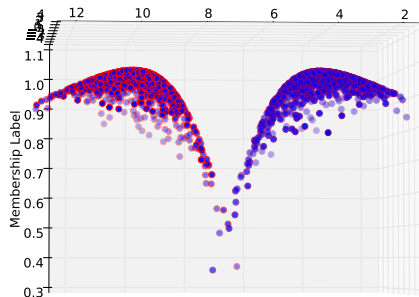
Here the clustering of two Gaussian Clusters





# Example

Here the clustering of two Gaussian Clusters



# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
- **Pros and Cons of FCM**
- What can we do? Possibilistic Clustering
  - Cost Function

# Pros and Cons of Fuzzy C-Means

## Advantages

- Unsupervised
- Always converges

# Pros and Cons of Fuzzy C-Means

## Advantages

- Unsupervised
- Always converges

## Disadvantages

- Long computational time
- Sensitivity to the initial guess (speed, local minima)
- Sensitivity to noise
  - ▶ One expects low (or even no) membership degree for outliers (noisy points)

# Pros and Cons of Fuzzy C-Means

## Advantages

- Unsupervised
- Always converges

## Disadvantages

- Long computational time
- Sensitivity to the initial guess (speed, local minima)
- Sensitivity to noise
  - ▶ One expects low (or even no) membership degree for outliers (noisy points)

# Pros and Cons of Fuzzy C-Means

## Advantages

- Unsupervised
- Always converges

## Disadvantages

- Long computational time
- Sensitivity to the initial guess (speed, local minima)
- Sensitivity to noise
  - ▶ One expects low (or even no) membership degree for outliers (noisy points)

# Pros and Cons of Fuzzy C-Means

## Advantages

- Unsupervised
- Always converges

## Disadvantages

- Long computational time
- Sensitivity to the initial guess (speed, local minima)
- Sensitivity to noise

► One expects low (or even no) membership degree for outliers (noisy points)

# Pros and Cons of Fuzzy C-Means

## Advantages

- Unsupervised
- Always converges

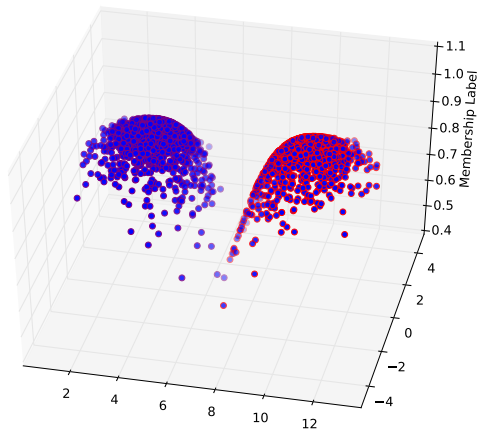
## Disadvantages

- Long computational time
- Sensitivity to the initial guess (speed, local minima)
- Sensitivity to noise
  - ▶ One expects low (or even no) membership degree for outliers (noisy points)



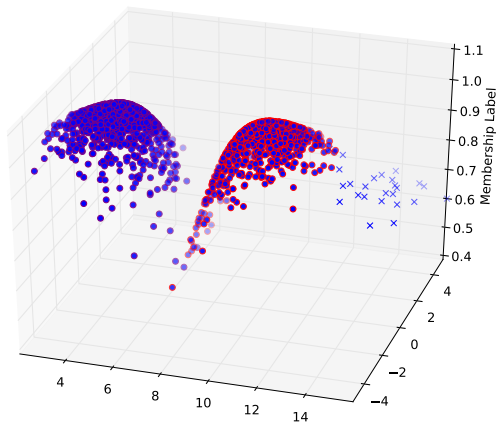
# Outliers, Disadvantage of FCM

After running without outliers



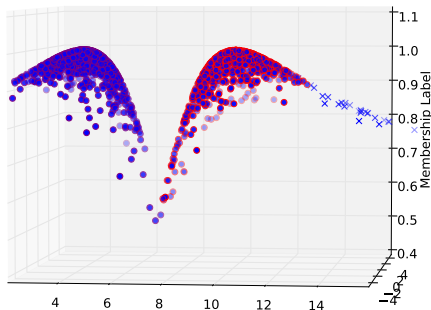
# Outliers, Disadvantage of FCM

Now add outliers (Shown in blue  $x$ 's) and their high memberships



# Outliers, Disadvantage of FCM

Now add outliers (Shown in blue  $x'$ 's) and their high memberships



# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

## Following Zadeh

They took in consideration that each class prototype as defining an elastic constraint.

What?

Giving the  $t_i(x_k)$  as degree of compatibility of sample  $x_k$  with cluster  $C_i$ .

We do the following:

If we consider the  $C_i$  as fuzzy sets over the set of samples

$X = \{x_1, x_2, \dots, x_N\}$

## Following Zadeh

They took in consideration that each class prototype as defining an elastic constraint.

## What?

Giving the  $t_i(\mathbf{x}_k)$  as degree of compatibility of sample  $\mathbf{x}_k$  with cluster  $C_i$ .

What's the following:

If we consider the  $C_i$  as fuzzy sets over the set of samples

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

## Following Zadeh

They took in consideration that each class prototype as defining an elastic constraint.

## What?

Giving the  $t_i(\mathbf{x}_k)$  as degree of compatibility of sample  $\mathbf{x}_k$  with cluster  $C_i$ .

## We do the following

If we consider the  $C_i$  as fuzzy sets over the set of samples

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

## Here is the Catch!!!

We should not use the old membership

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = 1 \quad (45)$$

Because

This is quite probabilistic... which is not what we want!!!

Thus

We only ask for membership, now using the possibilistic notation of  $t_i(\mathbf{x}_k)$ . (This is known as **typicality** value), to be in the interval  $[0, 1]$ .



## Here is the Catch!!!

We should not use the old membership

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = 1 \quad (45)$$

Because

This is quite probabilistic... which is not what we want!!!

This

We only ask for membership, now using the possibilistic notation of  $t_i(\mathbf{x}_k)$  (This is known as **typicality** value), to be in the interval  $[0, 1]$ .

## Here is the Catch!!!

We should not use the old membership

$$\sum_{i=1}^C A_i(\mathbf{x}_k) = 1 \quad (45)$$

Because

This is quite probabilistic... which is not what we want!!!

Thus

We only ask for membership, now using the possibilistic notation of  $t_i(\mathbf{x}_k)$  (This is known as **typicality** value), to be in the interval  $[0, 1]$ .

# New Constraints

First

$$t_i(\mathbf{x}_k) \in [0, 1] \quad \forall i, k \quad (46)$$

Second

$$0 < \sum_{k=1}^N t_i(\mathbf{x}_k) < N \quad \forall i \quad (47)$$

Third

$$\max_i t_i(\mathbf{x}_k) > 0 \quad \forall k \quad (48)$$

# New Constraints

First

$$t_i(\mathbf{x}_k) \in [0, 1] \quad \forall i, k \quad (46)$$

Second

$$0 < \sum_{k=1}^N t_i(\mathbf{x}_k) < N \quad \forall i \quad (47)$$

Third

$$\max_i t_i(\mathbf{x}_k) > 0 \quad \forall k \quad (48)$$

# New Constraints

First

$$t_i(\mathbf{x}_k) \in [0, 1] \quad \forall i, k \quad (46)$$

Second

$$0 < \sum_{k=1}^N t_i(\mathbf{x}_k) < N \quad \forall i \quad (47)$$

Third

$$\max_i t_i(\mathbf{x}_k) > 0 \quad \forall k \quad (48)$$

# Outline

## 1 $K$ -Means Clustering

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

## 2 $K$ -Meoids

- Introduction
- The Algorithm
- Complexity

## 3 The $K$ -Center Criterion Clustering

- Introduction
- Re-Stating the  $K$ -center as a Clustering Problem
- Comparison with  $K$ -means
- The Greedy  $K$ -Center Algorithm
- Pseudo-Code
- The  $K$ -Center Algorithm
- Notes in Implementation
- Examples
- $K$ -Center Algorithm Properties
- $K$ -Center Algorithm proof of correctness

## 4 Variations

- Fuzzy Clustering
  - Rethinking  $K$ -Means Cost Function
  - Using the Lagrange Multipliers
  - Examples
  - Pros and Cons of FCM
- What can we do? Possibilistic Clustering
  - Cost Function

We have the following cost function

### Cost Function

$$\sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (49)$$

### Problem

Unconstrained optimization of first term will lead to the trivial solution  $t_i(\mathbf{x}_k) = 0$  for all  $i, k$ .

Thus, we can introduce the following constraint

$$t_i(\mathbf{x}_k) \rightarrow 1 \quad (50)$$

- Roughly it means to make the typicality values as large as possible.

We have the following cost function

### Cost Function

$$\sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (49)$$

### Problem

Unconstrained optimization of first term will lead to the trivial solution  $t_i(\mathbf{x}_k) = 0$  for all  $i, k$ .

Thus, we can introduce the following constraint

$$t_i(\mathbf{x}_k) \rightarrow 1 \quad (50)$$

- Roughly it means to make the typicality values as large as possible.



We have the following cost function

### Cost Function

$$\sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (49)$$

### Problem

Unconstrained optimization of first term will lead to the trivial solution  $t_i(\mathbf{x}_k) = 0$  for all  $i, k$ .

Thus, we can introduce the following constraint

$$t_i(\mathbf{x}_k) \rightarrow 1 \quad (50)$$

- Roughly it means to make the typicality values as large as possible.

We can try to control this tendency

By putting all them together in

$$\sum_{k=1}^N (1 - t_i(\mathbf{x}_k))^m \quad (51)$$

With  $m$  to control the tendency of  $t_i(\mathbf{x}_k) \rightarrow 1$

We can also run this tendency over all the cluster using a suitable  $w_i > 0$  per cluster

$$\sum_{i=1}^C w_i \sum_{k=1}^N (1 - t_i(\mathbf{x}_k))^m \quad (52)$$

We can try to control this tendency

By putting all them together in

$$\sum_{k=1}^N (1 - t_i(\mathbf{x}_k))^m \quad (51)$$

With  $m$  to control the tendency of  $t_i(\mathbf{x}_k) \rightarrow 1$

We can also run this tendency over all the cluster using a suitable  $w_i > 0$  per cluster

$$\sum_{i=1}^C w_i \sum_{k=1}^N (1 - t_i(\mathbf{x}_k))^m \quad (52)$$

# Possibilistic C-Mean Clustering (PCM)

## The final Cost Function

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 + \sum_{i=1}^C w_i \sum_{k=1}^N (1 - t_i(\mathbf{x}_k))^m \quad (53)$$

Where

- $t_i(\mathbf{x}_k)$  are typicality values.
- $w_i$  are cluster weights

# Possibilistic C-Mean Clustering (PCM)

## The final Cost Function

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 + \sum_{i=1}^C w_i \sum_{k=1}^N (1 - t_i(\mathbf{x}_k))^m \quad (53)$$

## Where

- $t_i(\mathbf{x}_k)$  are **typicality** values.
- $w_i$  are cluster weights

# Explanation

## First Term

$$\sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (54)$$

It demands that the distance from feature vector to prototypes be as small as possible!!!

# Explanation

## First Term

$$\sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (54)$$

It demands that the distance from feature vector to prototypes be as small as possible!!!

## Second Term

$$\sum_{i=1}^c w_i \sum_{k=1}^n (1 - t_i(x_k))^m \quad (55)$$

It forces the typicality values  $t_i(x_k)$  to be as large as possible.

# Explanation

## First Term

$$\sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (54)$$

It demands that the distance from feature vector to prototypes be as small as possible!!!

## Second Term

$$\sum_{i=1}^c w_i \sum_{k=1}^n (1 - t_i(x_k))^m \quad (55)$$

It forces the typicality values  $t_i(x_k)$  to be as large as possible.



# Explanation

## First Term

$$\sum_{k=1}^N \sum_{i=1}^C [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (54)$$

It demands that the distance from feature vector to prototypes be as small as possible!!!

## Second Term

$$\sum_{i=1}^c w_i \sum_{k=1}^n (1 - t_i(x_k))^m \quad (55)$$

It forces the typicality values  $t_i(x_k)$  to be as large as possible.

# Final Updating Equations

## Typicality Values

$$t_i(\mathbf{x}_k) = \frac{1}{1 + \left(\frac{\|\mathbf{x}_k - \mathbf{v}_i\|^2}{w_i}\right)^{\frac{1}{m-1}}}, \quad \forall i, k \quad (56)$$

## Cluster Centers

$$\mathbf{v}_i = \frac{\sum_{k=1}^N t_i(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N t_i(\mathbf{x}_k)^m} \quad (57)$$

# Final Updating Equations

## Typicality Values

$$t_i(\mathbf{x}_k) = \frac{1}{1 + \left(\frac{\|\mathbf{x}_k - \mathbf{v}_i\|^2}{w_i}\right)^{\frac{1}{m-1}}}, \quad \forall i, k \quad (56)$$

## Cluster Centers

$$\mathbf{v}_i = \frac{\sum_{k=1}^N t_i(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^n t_i(\mathbf{x}_k)^m} \quad (57)$$

# Final Updating Equations

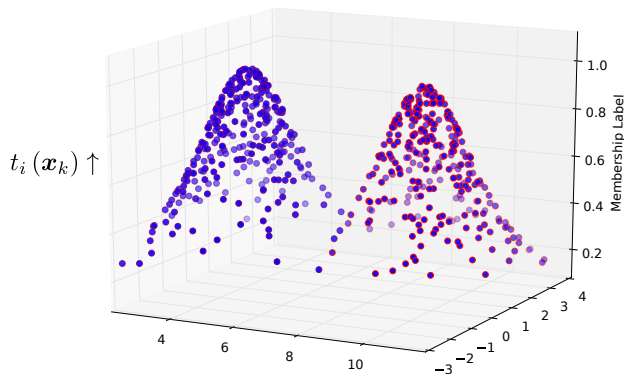
## Weights

$$w_i = M \frac{\sum_{k=1}^N [t_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2}{\sum_{k=1}^n [t_i(\mathbf{x}_k)]^m}, \quad (58)$$

with  $M > 0$ .

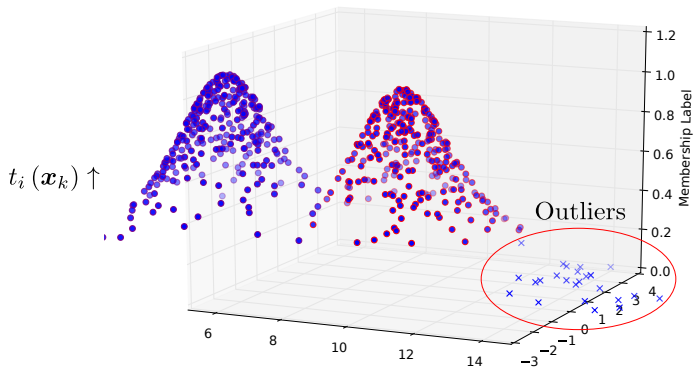
## Possibilistic can deal with outliers

Two Gaussian Clusters with  $\mu_1 = (4, 0)^T$ ,  $\mu_2 = (10, 0)^T$  and  $\sigma^2 = 1.0$



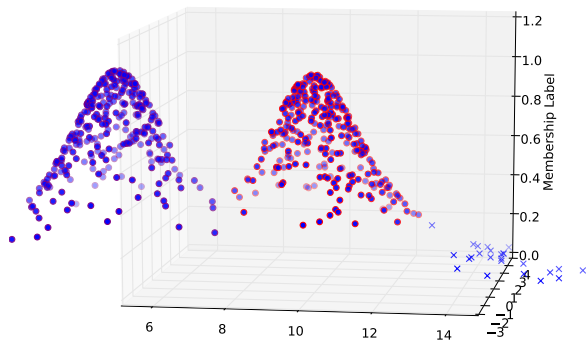
# Possibilistic can deal with outliers

Now add outliers



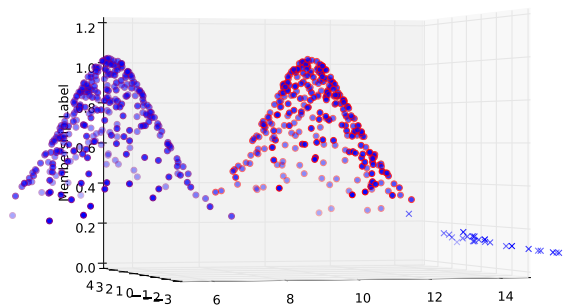
# Possibilistic can deal with outliers

## Another Angle



# Possibilistic can deal with outliers

## Another Angle





# Pros and Cons of Fuzzy $C$ -Means

## Advantages

Clustering noisy data samples.

# Pros and Cons of Fuzzy $C$ -Means

## Advantages

Clustering noisy data samples.

## Disadvantages

- Very sensitive to good initialization.
- For example, I needed to run fuzzy  $C$ -means to obtain good initial typicalities

# Pros and Cons of Fuzzy $C$ -Means

## Advantages

Clustering noisy data samples.

## Disadvantages

- Very sensitive to good initialization.
- For example, I needed to run fuzzy  $C$ -means to obtain good initial typicalities

## In Between!!!

Coincident clusters may result.

- Because the columns and rows of the typicality matrix are independent of each other.
- This could be advantageous (start with a large value of  $C$  and get less distinct clusters)

# Pros and Cons of Fuzzy $C$ -Means

## Advantages

Clustering noisy data samples.

## Disadvantages

- Very sensitive to good initialization.
- For example, I needed to run fuzzy  $C$ -means to obtain good initial typicalities

## In Between!!!

Coincident clusters may result.

- Because the columns and rows of the typicality matrix are independent of each other.
- This could be advantageous (start with a large value of  $C$  and get less distinct clusters)

# Pros and Cons of Fuzzy $C$ -Means

## Advantages

Clustering noisy data samples.

## Disadvantages

- Very sensitive to good initialization.
- For example, I needed to run fuzzy  $C$ -means to obtain good initial typicalities

## In Between!!!

Coincident clusters may result.

- Because the columns and rows of the typicality matrix are independent of each other.
- This could be advantageous (start with a large value of  $C$  and get less distinct clusters)

# Pros and Cons of Fuzzy $C$ -Means

## Advantages

Clustering noisy data samples.

## Disadvantages

- Very sensitive to good initialization.
- For example, I needed to run fuzzy  $C$ -means to obtain good initial typicalities

## In Between!!!

Coincident clusters may result.

- Because the columns and rows of the typicality matrix are independent of each other.
- This could be advantageous (start with a large value of  $C$  and get less distinct clusters)

## Nevertheless

There are more advanced clustering methods based on the possibilistic and fuzzy idea

Pal, N.R.; Pal, K.; Keller, J.M.; Bezdek, J.C., "A Possibilistic Fuzzy c-Means Clustering Algorithm," Fuzzy Systems, IEEE Transactions on , vol.13, no.4, pp.517,530, Aug. 2005.