

# Clustering algorithms

Konstantinos Koutroumbas

## Unit 4

- CFO clustering algorithms (k-means)

# CFO clustering algorithms: A unified view

## Data

$$X = \{\mathbf{x}_j \in R^l, j = 1, \dots, N\}$$

## Basic parameters - notation

- ✓  $\Theta = \{\boldsymbol{\theta}_j, j = 1, \dots, m\}$  ( $\boldsymbol{\theta}_j$  is the **representative** of cluster  $C_j$ ).
  - **Proximity** between  $\mathbf{x}_i$  and  $C_j$ :  $d(\mathbf{x}_i, \boldsymbol{\theta}_j)$

# CFO clustering algorithms: A unified view

## Basic parameters – notation (cont.)

$$\checkmark \quad U = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1m} \\ u_{21} & u_{22} & \cdots & u_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ u_{N1} & u_{N2} & \cdots & u_{Nm} \end{bmatrix} \equiv \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_N \end{bmatrix}$$

In the **probabilistic** case  
 $u_{ij}$  stands for  $P(j|\mathbf{x}_i)$

- $u_{ij} \in [0,1]$  quantifies the “**relation**” between  $\mathbf{x}_i$  and  $C_j$ .
- “**Large**” (“**small**”)  $u_{ij}$  values indicate **close** (**loose**) **relation** between  $\mathbf{x}_i$  and  $C_j$ .

$\Rightarrow u_{ij}$  varies **inversely proportional** wrt  $d(\mathbf{x}_i, \theta_j)$ .

- $\mathbf{u}_i$  : vector containing the  $u_{ij}$ 's of  $\mathbf{x}_i$  with all clusters.

-----

(\*) Unless otherwise stated, the case where **cluster representatives** are used is considered.

# CFO clustering algorithms: A unified view

## Aim:

- ✓ To **place** the **representatives** into dense in data regions (**physical clusters**).

## How this is achieved:

- ✓ Via the **minimization** of the following type of cost function (wrt  $\Theta, U$ )

$$J(\Theta, U) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(\mathbf{x}_i, \boldsymbol{\theta}_j) \quad (q \geq 1)$$

s.t. some **constraints** on  $U, C(U)$ .

For the **probabilistic** case  $d(\mathbf{x}_i, \boldsymbol{\theta}_j)$  is embedded in the **log-likelihood** of suitably defined **exponential distributions**

## Intuition:

- ✓ For **fixed**  $\boldsymbol{\theta}_j$ 's,  $J(\Theta, U)$  is a weighted sum of **fixed** distances  $d(\mathbf{x}_i, \boldsymbol{\theta}_j)$ .
- ⇒ **Minimization** of  $J(\Theta, U)$  wrt  $u_{ij}$  instructs for **large** weights ( $u_{ij}$ ) for **small** distances  $d(\mathbf{x}_i, \boldsymbol{\theta}_j)$ .
- ✓ For **fixed**  $u_{ij}$ 's, **minimization** of  $J(\Theta, U)$  wrt  $\boldsymbol{\theta}_j$ 's leads  $\boldsymbol{\theta}_j$ 's closer to their most relative data points.

# CFO clustering algorithms: A unified view

Basic types of algorithms:

**Constraints on  $U = [u_{ij}]$**

*Partition matrix*

*Membership matrix*

*Compatibility matrix*

**Hard:**

- $u_{ij} \in \{0, 1\}$
- $\sum_{j=1}^m u_{ij} = 1$

**Fuzzy:**

- $u_{ij} \in (0, 1)$
- $\sum_{j=1}^m u_{ij} = 1$

**Possibilistic (>1 choices):**

- $u_{ij} \in (0, 1]$

k-means

FCV

FCL

FOM

PCM

APCH



*k-dim. nonlinear manifold*

*k-dim. lin. manifold*

*Compact set in k-dim. lin. manifold*

$\Theta = \{\theta_j, j = 1, \dots, m\}$

# CFO clustering algorithms: A unified view

“Array of CFO algorithms”

$C(U)$

algorithm

$\theta_j$

	Hard Constr.	Fuzzy Constr.	Possib. Constr.	...
Point				
Line				
Hyperplane				
Hyperellipsoid				
...				

There are **several unexplored areas** (groups of algorithms) in this array.

# CFO clustering algorithms: A unified view

## General cost function opt. (CFO) scheme:

- ✓ Initialize  $\Theta = \Theta(0)$
  
- ✓  $t = 0$
  
- ✓ **Repeat**
  - $U(t) = \operatorname{argmin}_U J(\Theta(t), U)$ , s.t.  $C(U(t))$
  
  - $t = t + 1$
  
  - $\Theta(t) = \operatorname{argmin}_\Theta J(\Theta, U(t - 1))$
  
- ✓ **Until convergence**

# CFO clustering algorithms: A unified view

“Array of CFO algorithms”

$C(U)$

	Hard Constr.	Fuzzy Constr.	Possib. Constr.	...
Point	Hard CFO scheme	Fuzzy CFO scheme	Possib. CFO scheme	
Line				
Hyperplane				
Hyperellipsoid				
...				



# CFO clustering algorithms: A unified view

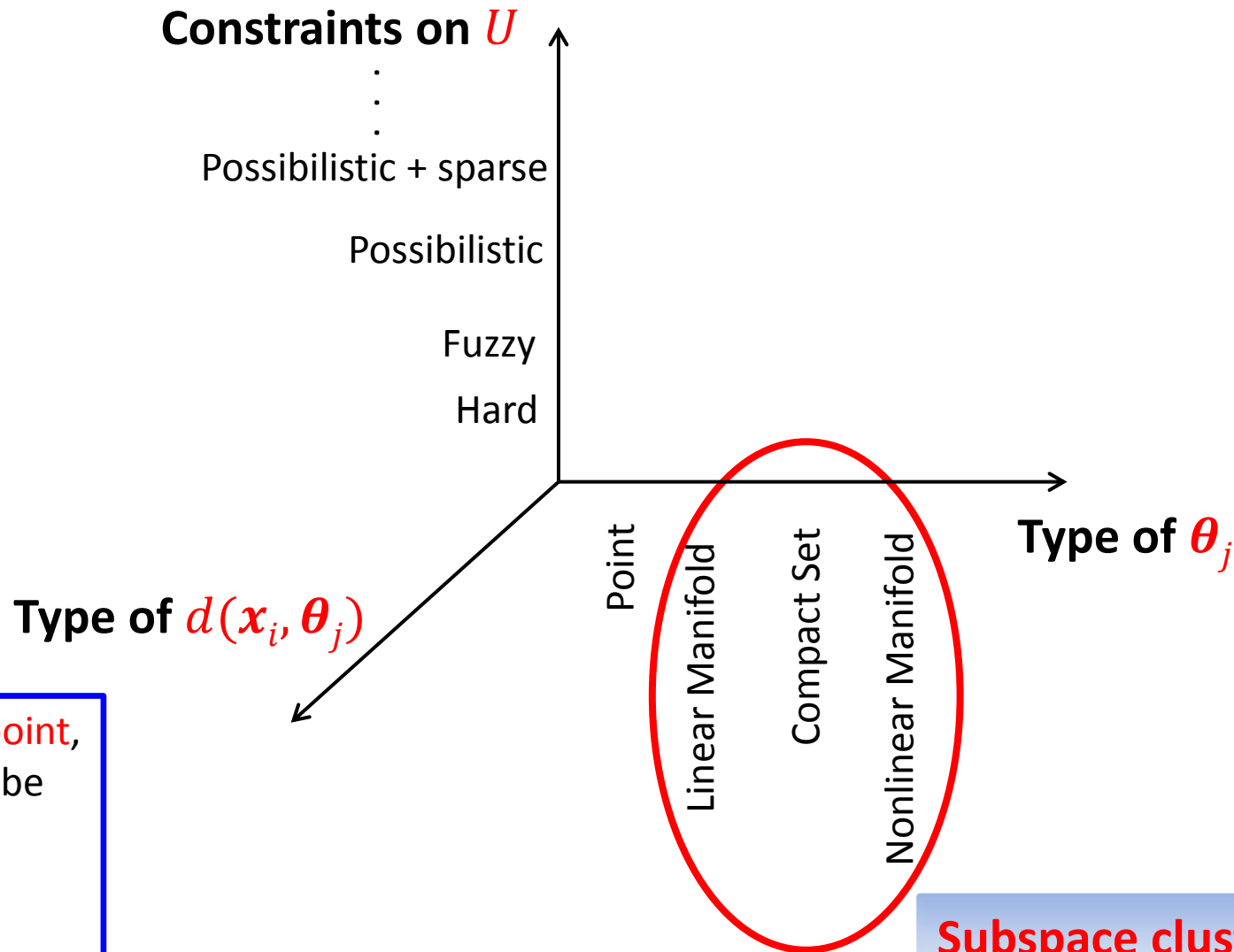
“Array of CFO algorithms”

$C(U)$

	Hard Constr.	Fuzzy Constr.	Possib. Constr.	...
Point	c-means scheme			
Line	c-lines scheme			
Hyperplane	c-hyperplanes scheme			
Hyperellipsoid	c-hyperellipsoids scheme			
...				

# CFO clustering algorithms: A unified view

## CFO clustering algorithms: A loose presentation



E.g.: If  $\theta_j$  is a point,  $d(x_i, \theta_j)$  may be

- Sq. Euclidean
- $l_p$  norm
- Mahalanobis

Subspace clustering

# CFO clustering algorithms: A unified view

“Array of CFO algorithms”

		$C(U)$			
		Hard Constr.	Fuzzy Constr.	Possib. Constr.	...
$\theta_j$	Point				
	Line				
	Hyperplane				
	Hyperellipsoid				
	...				

# Cost function optimization (CFO) algorithms

## Hard clustering algorithms:

Let  $X = \{x_1, x_2, \dots, x_N\}$  be a set of data points.

Each vector belongs **exclusively** to a single cluster.

Each **cluster** is **represented** by a representative  $\theta_j$  (point repr., hyperplane...).

Let  $\Theta = \{\theta_1, \theta_2, \dots, \theta_m\}$

Define  $u_{ij} = \begin{cases} 1, & \text{if } x_i \in C_j \\ 0, & \text{otherwise} \end{cases}$  and  $U = [u_{ij}]_{N \times m}$

It is  $\sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$

Define the **cost function**

$$J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} d(x_i, \theta_j) = \sum_{j=1}^m \sum_{x_i \in C_j} d(x_i, \theta_j)$$

When  $J(U, \Theta)$  is **minimized**?

# CFO hard clustering algorithms

$$J(U, \theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} d(\mathbf{x}_i, \theta_j) = \sum_{j=1}^m \sum_{\mathbf{x}_i \in C_j} d(\mathbf{x}_i, \theta_j)$$

For **fixed  $\theta_j$ 's**: When, for each  $\mathbf{x}_i$ , only its distance from its closest representative is taken into account.

This suggests to **define**  $u_{ij} = \begin{cases} 1, & \text{if } d(\mathbf{x}_i, \theta_j) = \min_{q=1, \dots, m} d(\mathbf{x}_i, \theta_q) \\ 0, & \text{otherwise} \end{cases}$

For **fixed  $u_{ij}$ 's**: Solve the following  **$m$**  independent problems

$$\min_{\theta_j} \sum_{\mathbf{x}_i \in C_j} d(\mathbf{x}_i, \theta_j) \equiv \min_{\theta_j} \sum_{i=1}^N u_{ij} d(\mathbf{x}_i, \theta_j)$$

Thus, the **Generalized Hard Algorithmic Scheme (GHAS)** is given below

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

- Choose  $\theta_j(0)$  as **initial estimates** for  $\theta_j, j=1, \dots, m$ .
- $t = 0$
- **Repeat**

– For  $i = 1$  to  $N$  % *Determination of the partition*

o For  $j = 1$  to  $m$

$$u_{ij}(t) = \begin{cases} 1, & \text{if } d(\mathbf{x}_i, \theta_j(t)) = \min_{q=1, \dots, m} d(\mathbf{x}_i, \theta_q(t)) \\ 0, & \text{otherwise} \end{cases}$$

o End {For- $j$ }

– End {For- $i$ }

–  $t = t + 1$

– For  $j = 1$  to  $m$  % *Parameter updating*

o Set

$$\theta_j(t) = \operatorname{argmin}_{\theta_j} \sum_{i=1}^N u_{ij}(t-1) d(\mathbf{x}_i, \theta_j), j = 1, \dots, m$$

– End {For- $j$ }

- **Until** a **termination criterion** is met.

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### Remarks:

- In the update of each  $\theta_j$ , only the vectors  $x_i$  for which  $u_{ij}(t - 1) = 1$  are used.
- **GHAS** may **terminate** when either
  - $\|\theta(t) - \theta(t - 1)\| < \varepsilon$  or
  - $U$  remains **unchanged** for **two successive iterations**.
- The two-step optimization procedure in GHAS **does not necessarily lead to a local minimum** of  $J(U, \theta)$ .

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

#### General comments

- It is a special case of GHAS where
  - **Point representatives** are **used**.
  - The **squared Euclidean distance** is **employed**.
- The cost function  $J(U, \Theta)$  becomes now
$$J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} \|\mathbf{x}_i - \boldsymbol{\theta}_j\|^2$$
- Applying GHAS in this case, it turns out that it **converges** to a **minimum** of the **cost function**.
- Isodata **recovers clusters** that are as **compact** as possible.
- For other choices of the distance (including the Euclidean), the algorithm converges but not necessarily to a minimum of  $J(U, \Theta)$ .



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

- Choose arbitrary initial estimates  $\theta_j(0)$  for the  $\theta_j$ 's,  $j=1, \dots, m$ .

- $t = 0$

- **Repeat**

- For  $i = 1$  to  $N$  % *Determination of the partition*

- o For  $j=1$  to  $m$

$$u_{ij}(t) = \begin{cases} 1, & \text{if } \|\mathbf{x}_i - \theta_j(t)\|^2 = \min_{q=1, \dots, m} \|\mathbf{x}_i - \theta_q(t)\|^2 \\ 0, & \text{otherwise} \end{cases}$$

- o End {For- $j$ }

- End {For- $i$ }

- $t = t + 1$

- For  $j = 1$  to  $m$  % *Parameter updating*

- o Set

$$\theta_j(t) = \frac{\sum_{i=1}^N u_{ij}(t-1) \mathbf{x}_i}{\sum_{i=1}^N u_{ij}(t-1)}, j = 1, \dots, m$$

- End {For- $j$ }

- **Until** no change in  $\theta_j$ 's occurs between two successive iterations

# CFO hard clustering algorithms

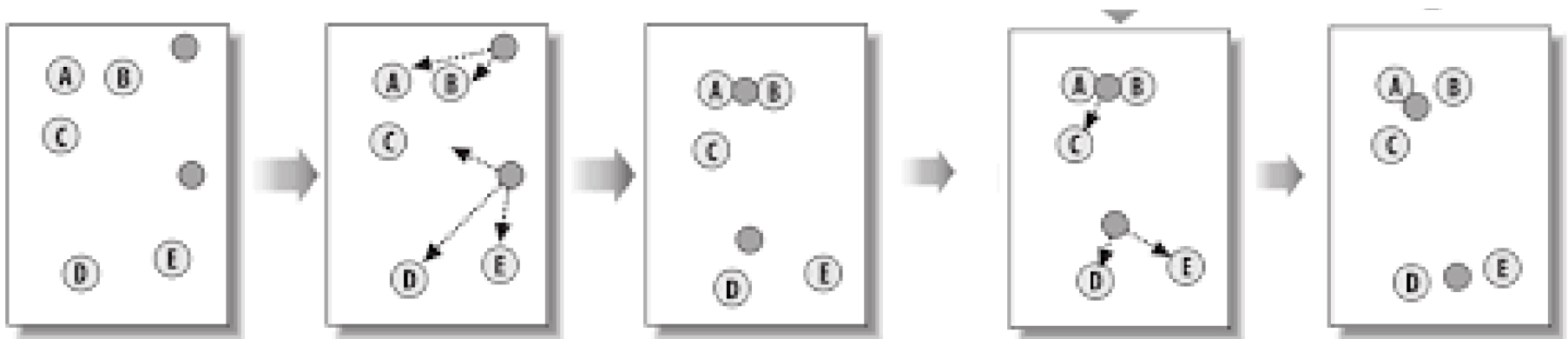
## The k-means case.

Choose arbitrary initial estimates  $\theta_j(0)$  for the  $\theta_j$ 's,  $j = 1, \dots, m$ .

### Repeat

- For  $i = 1$  to  $N$  *Partition determination*
  - o Determine the closest representative, say  $\theta_j$ , for  $x_i$
  - o Set  $u_{ij} = 1$  and  $u_{iq} = 0$ ,  $q = 1, \dots, m$ ,  $q \neq j$ .
- End {For}
- For  $j = 1$  to  $m$  *Parameter updating*
  - o Determine  $\theta_j$  as the mean of the vectors  $x_i \in X$  with  $u_{ij} = 1$ .
- End {For}

Until no change in  $\theta_j$ 's occurs between two successive iterations



# CFO hard clustering algorithms

## Remarks

- It is a **batch, single clustering** algorithm
- It is a **hard clustering** algorithm that uses **point representatives**  $\theta_j$  for the clusters  $C_j$ .
- It results from the optimization of the following cost function

$$J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} \|\mathbf{x}_i - \theta_j\|^2$$

where  $U = [u_{ij}]$  and  $\Theta = \{\theta_1, \dots, \theta_m\}$

- It is of **iterative** nature.
- **Initially** it places the representatives  $\theta_j$  at **random positions** in space.
- It gradually **moves the representatives** towards the **centers** of the **true clusters**.
- In practice, its **time complexity** is  $O(q \cdot m \cdot N)$  ( $q$  is the number of iterations).
- It requires the number of clusters  $m$  to be **known a priori**.

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

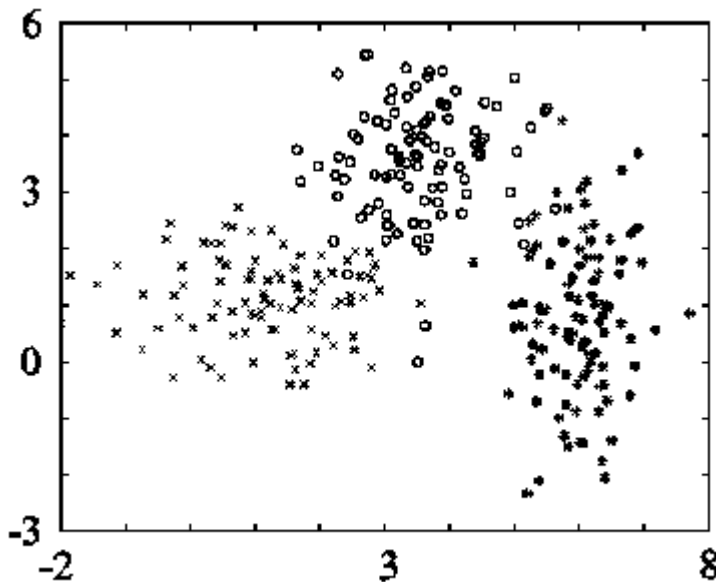
**Example 1:** (a) Consider three two-dimensional normal distributions with mean values:

$$\boldsymbol{\mu}_1 = [1, 1]^T, \boldsymbol{\mu}_2 = [3.5, 3.5]^T, \boldsymbol{\mu}_3 = [6, 1]^T$$

and respective covariance matrices

$$\Sigma_1 = \begin{bmatrix} 1 & -0.3 \\ -0.3 & 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0.3 \\ 0.3 & 1 \end{bmatrix}, \Sigma_3 = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}$$

Generate a group of **100 vectors** from **each distribution**. These form the data set  $X$ .



**Confusion matrix** for the results of k-means.

$$A = \begin{bmatrix} 94 & 3 & 3 \\ 0 & 100 & 0 \\ 9 & 0 & 91 \end{bmatrix}$$

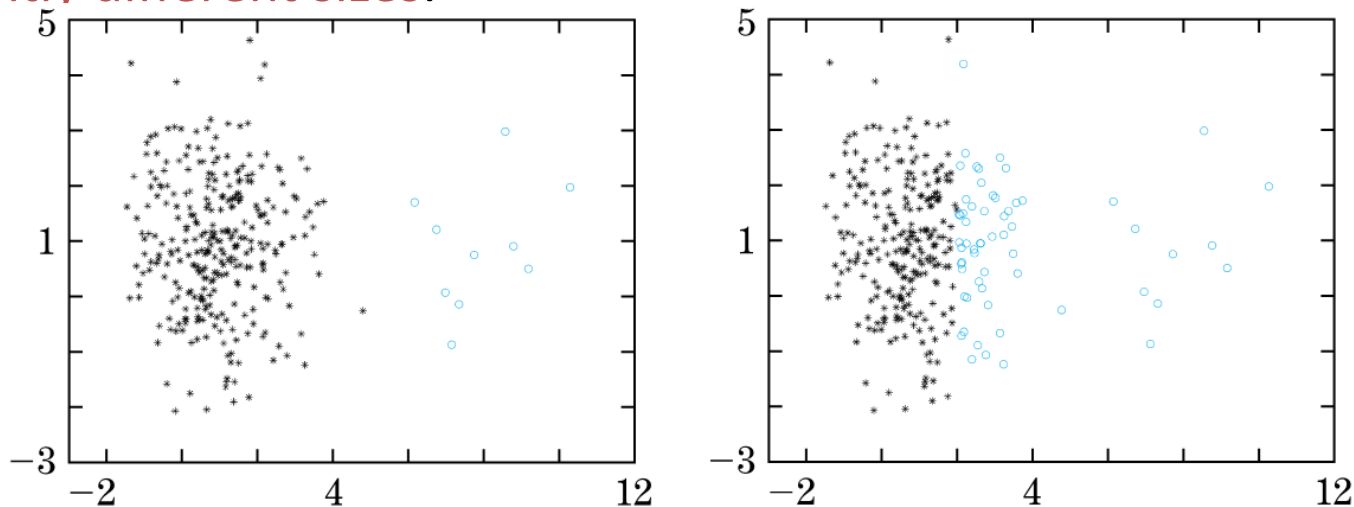
# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

**Example 2:** (i) Consider two 2-dimensional Gaussian distributions  $N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ ,  $N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ , with  $\boldsymbol{\mu}_1 = [1, 1]^T$ ,  $\boldsymbol{\mu}_2 = [8, 1]^T$ ,  $\boldsymbol{\Sigma}_1 = 1.5I$  and  $\boldsymbol{\Sigma}_2 = I$ . (ii) Generate **300 points** from the **1<sup>st</sup> distribution** and **10 points** from the **2<sup>nd</sup> distribution**. (iii) Set  $m = 2$  and initialize randomly  $\boldsymbol{\theta}_j$ 's ( $\boldsymbol{\theta}_j \equiv \boldsymbol{\mu}_j$ ).

- After convergence the large group has been split into two clusters.
- Its right part has been assigned to the same cluster with the points of the small group (see figure below).
- This indicates that **k-means cannot deal accurately with clusters having significantly different sizes**.



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

#### Remarks:

- $k$ -means recovers **compact clusters**.
- The computational complexity of the  $k$ -means is  $O(Nmq)$ , where  $q$  is the number of iterations required for convergence. In practice,  $m$  and  $q$  are significantly less than  $N$ , thus,  **$k$ -means becomes eligible for processing large data sets**.
- **Sequential (online) versions** of the  $k$ -means, where the updating of the representatives takes place immediately after the identification of the representative that lies closer to the current input vector  $\mathbf{x}_i$ , have also been proposed.
- A variant of the  $k$ -means results if the number of vectors in each cluster is constrained *a priori*.

#### Further remarks:

Some drawbacks of the original  $k$ -means accompanied with the variants of the  $k$ -means that deal with them are discussed next.

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

**Drawback 1:** *Different initial partitions may lead  $k$ -means to produce different final clusterings, each one corresponding to a different local minimum.*

### Strategies for facing drawback 1:

- Single run methods

- Use a sequential algorithm (discussed previously) to produce initial estimates for  $\theta_j$ 's.
- Partition randomly the data set into  $m$  subsets and use their means as initial estimates for  $\theta_j$ 's.

- Multiple run methods

- Create different partitions of  $X$ , run  $k$ -means for each one of them and select the best result.

- Utilization of tools from stochastic optimization techniques (simulated annealing, genetic algorithms etc).

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

**Drawback 2:** Knowledge of the number of clusters  $m$  is required a priori.

### Strategies for facing drawback 2:

- Employ splitting, merging and/or discarding operations of the clusters resulting from  $k$ -means.
- Estimate  $m$  as follows:
  - Run a **sequential** algorithm many times for different thresholds of dissimilarity  $\theta$ .
  - Plot  $\theta$  versus the number of clusters and identify the largest plateau in the graph and set  $m$  equal to the value that corresponds to this plateau.



# CFO hard clustering algorithms

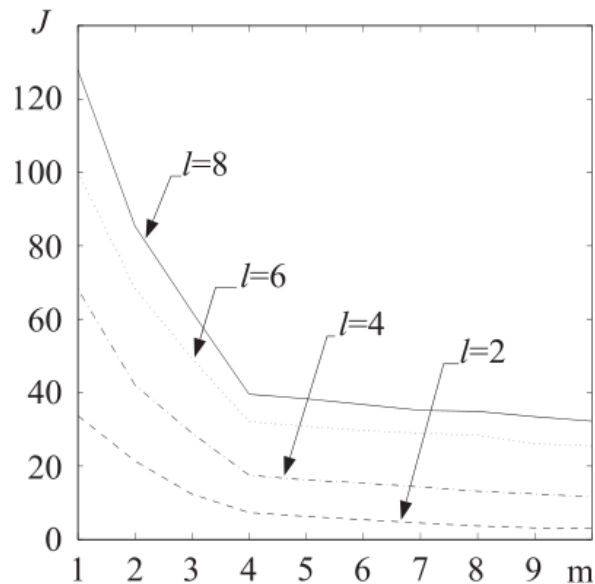
## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

**Drawback 2:** Knowledge of the number of clusters  $m$  is required a priori.

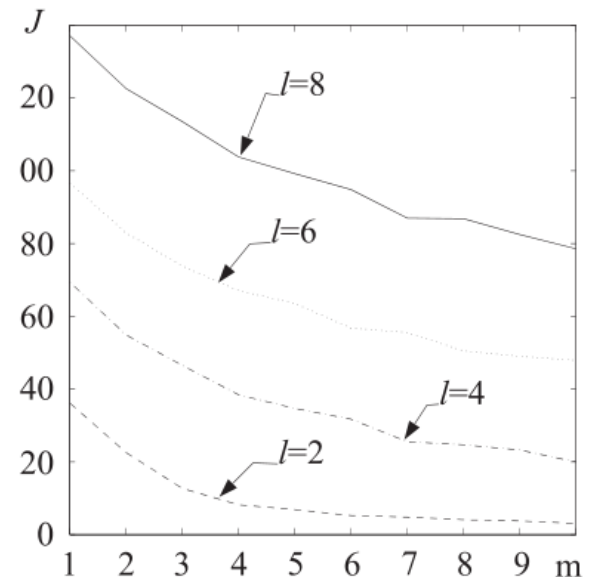
Strategies for facing drawback 2 (cont.):

- Estimate  $m$  as follows:
  - Run the  **$k$ -means** algorithm for different values of the number of clusters  $m$ .
  - For each of the resulting clusterings compute the value of  $J$ .
  - **Plot  $J$  versus** the number of clusters  $m$  and identify the most significant knee in the graph. Its position indicates the number of physical clusters.



Clustered data

Non-clustered data



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The Isodata or $k$ -Means or $c$ -Means algorithm

**Drawback 3:**  *$k$ -means is sensitive to outliers and noise.*

#### Strategies for facing drawback 3:

- Discard all “small” clusters (they are likely to be formed by outliers).
- Use a  $k$ -medoids algorithm (see below), where a cluster is represented by one of its points.

**Drawback 4:**  *$k$ -means is not suitable for data with nominal (categorical) coordinates.*

#### Strategies for facing drawback 4:

- Use a  $k$ -medoids algorithm.