

# Clustering algorithms

Konstantinos Koutroumbas

## Unit 5

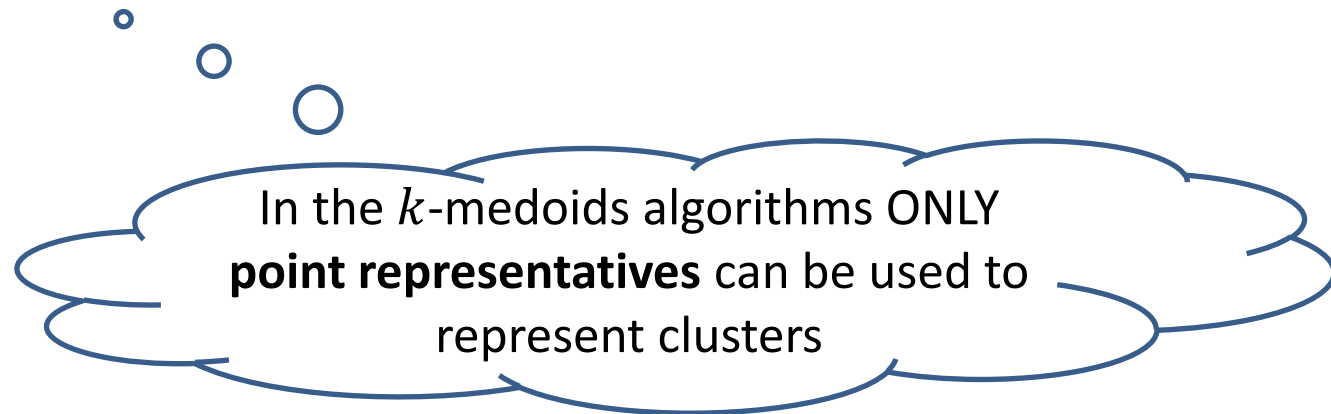
- Hard clustering: k-medoids algorithms (PAM, CLARA, CLARANS)
- Probabilistic clustering (EM algorithm)

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### *k*-Medoids Algorithms

- Each cluster is represented by a vector selected **among** the elements of  $X$  (**medoid**).
- A cluster contains
  - Its medoid
  - All vectors in  $X$  that
    - o Are not used as medoids in other clusters
    - o Lie closer to its medoid than the medoids representing other clusters.



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### *k-Medoids Algorithms*

Let

- $\Theta$  be the **set of medoids** of all clusters,
- $I_\Theta$  the set of **indices** of the points in  $X$  that constitute  $\Theta$  and
- $I_{X-\Theta}$  the set of indices of the points that are **not medoids**.

Obtaining the set of medoids  $\Theta$  that best represents the data set,  $X$  is equivalent to minimizing the following cost function

$$J(\Theta, U) = \sum_{i \in I_{X-\Theta}} \sum_{j \in I_\Theta} u_{ij} d(\mathbf{x}_i, \mathbf{x}_j)$$

with

$$u_{ij} = \begin{cases} 1, & \text{if } d(\mathbf{x}_i, \mathbf{x}_j) = \min_{q \in I_\Theta} d(\mathbf{x}_i, \mathbf{x}_q), \\ 0, & \text{otherwise} \end{cases}, \quad i = 1, \dots, N$$

# CFO hard clustering algorithms

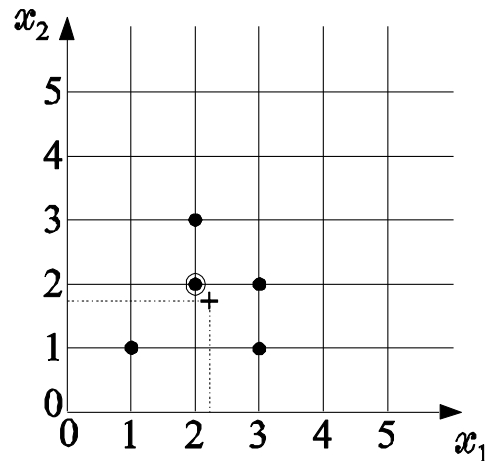
## Generalized Hard Algorithmic Scheme (GHAS)

### *k*-Medoids Algorithms

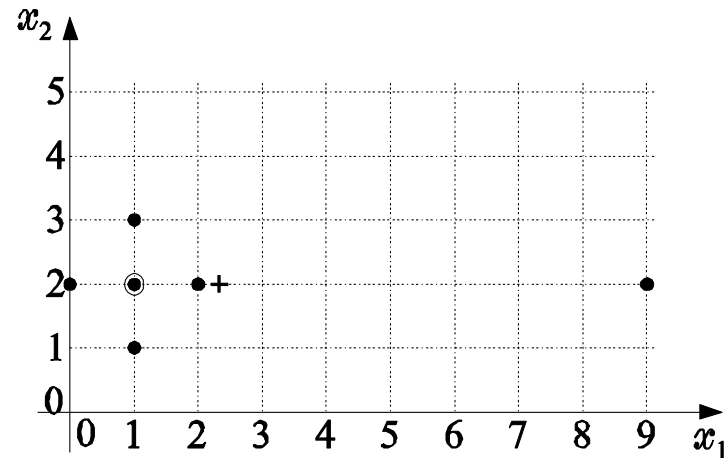
#### Example 3:

(a) The five-point two-dimensional set stems from the discrete domain  $D = \{1,2,3,4, \dots\} \times \{1,2,3,4, \dots\}$ . Its medoid is the circled point and **its mean** is the “+” point, which **does not belong to  $D$** .

(b) In the six-point two-dimensional set, the point (9,2) can be considered as an outlier. While **the outlier affects significantly the mean** of the set, **it does not affect its medoid**.



(a)



(b)

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

Representing clusters with **mean values** vs representing clusters with **medoids**

Mean Values	Medoids
1. Suited only for continuous domains	<b>1. Suited for either cont. or discrete domains</b>
2. Algorithms using means are sensitive to outliers	<b>2. Algorithms using medoids are less sensitive to outliers</b>
<b>3. The mean possess a clear geometrical and statistical meaning</b>	3. The medoid has not a clear geometrical meaning
<b>4. Algorithms using means are not computationally demanding</b>	4. Algorithms using medoids are more computationally demanding

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### *k*-Medoids Algorithms

#### Algorithms to be considered

- **PAM** (Partitioning Around Medoids)
- **CLARA** (Clustering LARge Applications)
- **CLARANS** (Clustering Large Applications based on RANdomized Search)

#### The PAM algorithm

- The number of clusters  $m$  is **required *a priori***.

## Definitions-preliminaries

- Two sets of medoids  $\Theta$  and  $\Theta'$ , each one consisting of  $m$  elements, are called **neighbors** if they **share**  $m - 1$  elements.
- A set  $\Theta$  of medoids with  $m$  elements can have  $m(N - m)$  neighbors.
- Let  $\Theta_{ij}$  denote the **neighbor** of  $\Theta$  that results if  $x_j, j \in I_{X-\Theta}$  **replaces**  $x_i, i \in I_{\Theta}$ .
- Let  $\Delta J_{ij} = J(\Theta_{ij}, U_{ij}) - J(\Theta, U)$ .

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm

- Determination of  $\Theta$  that best represents the data
  - **Generate** a set  $\Theta$  of  $m$  medoids, randomly selected out of  $X$ .
  - **(A) Determine** the neighbor  $\Theta_{qr}$ ,  $q \in I_{\Theta}$ ,  $r \in I_{X-\Theta}$  among the  $m(N - m)$  neighbors of  $\Theta$  for which  $\Delta J_{qr} = \min_{i \in I_{\Theta}, j \in I_{X-\Theta}} \Delta J_{ij}$ .
  - If  $\Delta J_{qr} < 0$  then
    - $\Delta J_{qr} < 0 \Leftrightarrow J(\Theta_{qr}, U_{qr}) - J(\Theta, U) < 0$
    - $\Leftrightarrow J(\Theta_{qr}, U_{qr}) < J(\Theta, U)$
    - Replace  $\Theta$  by  $\Theta_{qr}$
    - Go to **(A)**
  - End
- Assignment of points to clusters
  - Assign each  $x \in X - \Theta$  to the cluster represented by the closest to  $x$  medoid.

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm

Computation of  $\Delta J_{ij}$ .

It is defined as:

$$\begin{aligned}\Delta J_{ij} &= J(\Theta_{ij}, U_{ij}) - J(\Theta, U) = \sum_{s \in I_{X-\Theta_{ij}}} \sum_{t \in I_{\Theta_{ij}}} u_{st} d(\mathbf{x}_s, \mathbf{x}_t) - \sum_{s \in I_{X-\Theta}} \sum_{t \in I_{\Theta}} u_{st} d(\mathbf{x}_s, \mathbf{x}_t) \\ &\equiv \sum_{h \in I_{X-\Theta}} C_{hij}\end{aligned}$$

where  $C_{hij}$  is the difference in  $J$ , resulting from the (possible) assignment of the vector  $\mathbf{x}_h \in X - \Theta$  from the cluster it currently belongs to another, as a consequence of the replacement of  $\mathbf{x}_i \in \Theta$  by  $\mathbf{x}_j \in X - \Theta$ .

For the computation of  $C_{hij}$  associated with a specific  $\mathbf{x}_h \in X - \Theta$  it is required

- The **distance** of  $\mathbf{x}_h$  from its **closest medoid** in  $\Theta$
- The **distance** of  $\mathbf{x}_h$  from its **next to closest medoid** in  $\Theta$ .
- The **distance** of  $\mathbf{x}_h$  from the **newly inserted medoid** in  $\Theta_{ij}$ ,  $\mathbf{x}_j$ .



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm (cont.)

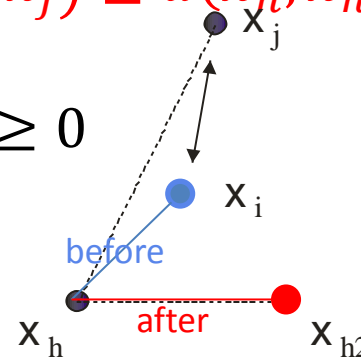
Computation of  $C_{hij}$ :

$x_h$  belongs to the cluster represented by  $x_i$  ( $x_{h2} \in \Theta$  denotes the second closest to  $x_h$  representative) and  $d(x_h, x_j) \geq d(x_h, x_{h2})$  ( $\geq d(x_h, x_i)$ ). Then

$$C_{hij} = d(x_h, x_{h2}) - d(x_h, x_i) \geq 0$$

Contribution of  $x_h$  to  $J(\Theta_{ij}, U_{ij})$

Contribution of  $x_h$  to  $J(\Theta, U)$

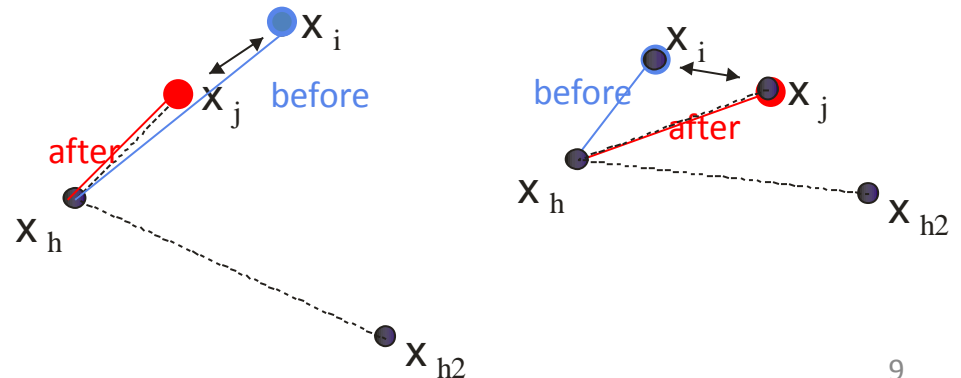


$x_h$  belongs to the cluster represented by  $x_i$  ( $x_{h2} \in \Theta$  denotes the second closest to  $x_h$  representative) and  $d(x_h, x_j) \leq d(x_h, x_{h2})$ . Then

$$C_{hij} = d(x_h, x_j) - d(x_h, x_i) (><) 0$$

Contribution of  $x_h$  to  $J(\Theta_{ij}, U_{ij})$

Contribution of  $x_h$  to  $J(\Theta, U)$



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm (cont.)

Computation of  $C_{hij}$  (cont.):

$x_h$  is not represented by  $x_i$  ( $x_{h1}$  denotes the closest to  $x_h$  medoid) and  $d(x_h, x_{h1}) \leq d(x_h, x_j)$ . Then

$$C_{hij} = d(x_h, x_{h1}) - d(x_h, x_{h1}) = 0$$

Contribution of  $x_h$  to  $J(\theta_{ij}, U_{ij})$

Contribution of  $x_h$  to  $J(\theta, U)$

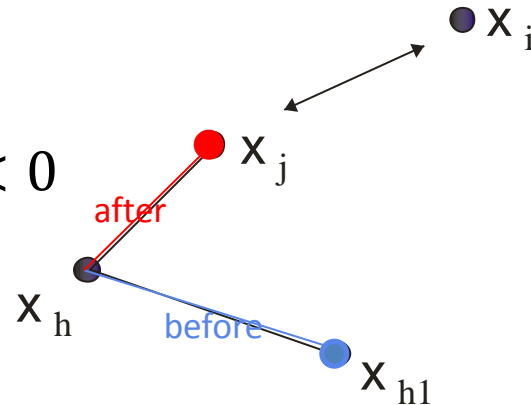


$x_h$  is not represented by  $x_i$  ( $x_{h1}$  denotes the closest to  $x_h$  medoid) and  $d(x_h, x_{h1}) > d(x_h, x_j)$ . Then

$$C_{hij} = d(x_h, x_j) - d(x_h, x_{h1}) < 0$$

Contribution of  $x_h$  to  $J(\theta_{ij}, U_{ij})$

Contribution of  $x_h$  to  $J(\theta, U)$



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm (cont.)

#### Remarks:

- Experimental results show the PAM works **satisfactorily with small data sets.**
- Its computational complexity is  $O(m(N - m)^2)$ . **Unsuitable for large data sets.**

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm (Example)

**Data set:**  $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ , with

$x_1 = [0,3]^T$ ,  $x_2 = [1,3]^T$ ,  $x_3 = [2,3]^T$ ,  $x_4 = [0,0]^T$ ,  $x_5 = [1,0]^T$ ,  $x_6 = [2,0]^T$ .

**Set of medoids:**  $\Theta = \{x_4, x_5\}$

Computation of  $J(\Theta, U)$  (Squared Euclidean distance is considered):

$$x_1 \rightarrow d(x_1, x_4) = 9 < 10 = d(x_1, x_5) \rightarrow u_{14} = 1, u_{15} = 0$$

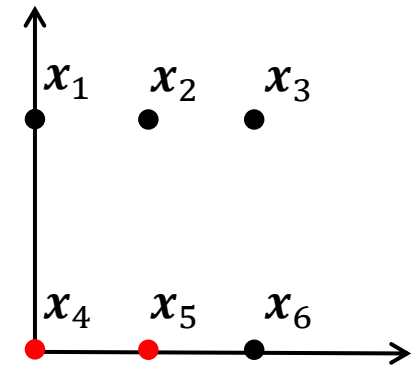
$$x_2 \rightarrow d(x_2, x_4) = 10 > 9 = d(x_2, x_5) \rightarrow u_{24} = 0, u_{25} = 1$$

$$x_3 \rightarrow d(x_3, x_4) = 13 > 10 = d(x_3, x_5) \rightarrow u_{34} = 0, u_{35} = 1$$

$$x_4 \rightarrow d(x_4, x_4) = 0 < 1 = d(x_4, x_5) \rightarrow u_{44} = 1, u_{45} = 0$$

$$x_5 \rightarrow d(x_5, x_4) = 1 > 0 = d(x_5, x_5) \rightarrow u_{54} = 0, u_{55} = 1$$

$$x_6 \rightarrow d(x_6, x_4) = 2 > 1 = d(x_6, x_5) \rightarrow u_{64} = 0, u_{65} = 1$$



$$\begin{aligned} J(\Theta, U) &= u_{14}d(x_1, x_4) + u_{15}d(x_1, x_5) + 1 \cdot 9 + 0 \cdot 10 + \\ &+ u_{24}d(x_2, x_4) + u_{25}d(x_2, x_5) + 0 \cdot 10 + 1 \cdot 9 + \\ &+ u_{34}d(x_3, x_4) + u_{35}d(x_3, x_5) + 0 \cdot 13 + 1 \cdot 10 + \\ &+ u_{44}d(x_4, x_4) + u_{45}d(x_4, x_5) + 1 \cdot 0 + 0 \cdot 1 + \\ &+ u_{54}d(x_5, x_4) + u_{55}d(x_5, x_5) + 0 \cdot 1 + 1 \cdot 0 + \\ &+ u_{64}d(x_6, x_4) + u_{65}d(x_6, x_5) + 0 \cdot 2 + 1 \cdot 1 \end{aligned} = 29$$

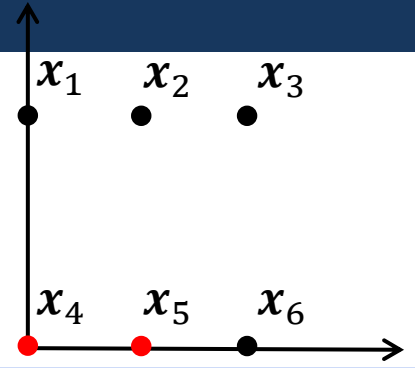
# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm (Example)

Data set:  $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ , with  
 $x_1 = [0,3]^T$ ,  $x_2 = [1,3]^T$ ,  $x_3 = [2,3]^T$ ,  $x_4 = [0,0]^T$ ,  $x_5 = [1,0]^T$ ,  $x_6 = [2,0]^T$ .

Set of medoids:  $\Theta = \{x_4, x_5\}$



$$\Theta_{42} = \{x_2, x_5\}$$

$$J(\Theta_{42}, U_{42}) = 4$$

$$\Delta J_{42} = 4 - 29 = -25$$

$$\Theta_{43} = \{x_3, x_5\}$$

$$J(\Theta_{43}, U_{43}) = 5$$

$$\Delta J_{43} = 5 - 29 = -24$$

$$\Theta_{46} = \{x_6, x_5\}$$

$$J(\Theta_{46}, U_{46}) = 29$$

$$\Delta J_{46} = 29 - 29 = 0$$

$$\Theta_{41} = \{x_1, x_5\}$$

$$J(\Theta_{41}, U_{41}) = 5$$

$$\Delta J_{41} = 5 - 29 = -24$$

$$\Theta = \{x_4, x_5\}$$

$$J(\Theta, U) = 29$$

$$\Theta_{51} = \{x_4, x_1\}$$

$$J(\Theta_{51}, U_{51}) = 6$$

$$\Delta J_{51} = 6 - 29 = -23$$

$$\Theta_{56} = \{x_4, x_6\}$$

$$J(\Theta_{56}, U_{56}) = 29$$

$$\Delta J_{56} = 29 - 29 = 0$$

$$\Theta_{53} = \{x_4, x_3\}$$

$$J(\Theta_{53}, U_{53}) = 5$$

$$\Delta J_{53} = 5 - 29 = -24$$

$$\Theta_{52} = \{x_4, x_2\}$$

$$J(\Theta_{52}, U_{52}) = 5$$

$$\Delta J_{52} = 5 - 29 = -24$$

It is  $\Delta J_{42} = \min_{i \in I_\Theta, j \in I_{X-\Theta}} \Delta J_{ij} = -25 < 0$   
 Thus, according to PAM,  $\Theta$  will be replaced by  $\Theta_{42}$ .

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm (Example)

**Data set:**  $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ , with

$x_1 = [0,3]^T$ ,  $x_2 = [1,3]^T$ ,  $x_3 = [2,3]^T$ ,  $x_4 = [0,0]^T$ ,  $x_5 = [1,0]^T$ ,  $x_6 = [2,0]^T$ .

**Set of medoids:**  $\Theta_{42} = \{x_2, x_5\}$

Computation of  $J(\Theta_{42}, U_{42})$  (Squared Euclidean distance is considered):

$$x_1 \rightarrow d(x_1, x_2) = 1 < 10 = d(x_1, x_5) \rightarrow u_{12} = 1, u_{15} = 0$$

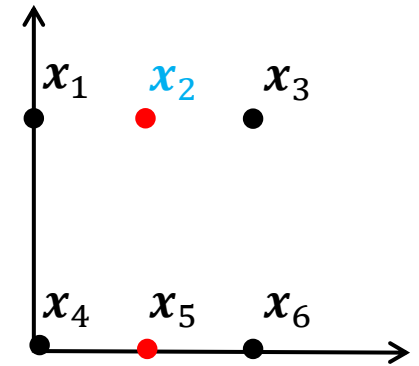
$$x_2 \rightarrow d(x_2, x_2) = 0 < 9 = d(x_2, x_5) \rightarrow u_{22} = 1, u_{25} = 0$$

$$x_3 \rightarrow d(x_3, x_2) = 1 < 10 = d(x_3, x_5) \rightarrow u_{32} = 1, u_{35} = 0$$

$$x_4 \rightarrow d(x_4, x_2) = 10 > 1 = d(x_4, x_5) \rightarrow u_{42} = 0, u_{45} = 1$$

$$x_5 \rightarrow d(x_5, x_2) = 9 > 0 = d(x_5, x_5) \rightarrow u_{52} = 0, u_{55} = 1$$

$$x_6 \rightarrow d(x_6, x_2) = 10 > 1 = d(x_6, x_5) \rightarrow u_{62} = 0, u_{65} = 1$$



$$\begin{aligned} J(\Theta_{42}, U_{42}) &= u_{12}d(x_1, x_2) + u_{15}d(x_1, x_5) + 1 \cdot 1 + 0 \cdot 10 + \\ &+ u_{22}d(x_2, x_2) + u_{25}d(x_2, x_5) + 1 \cdot 0 + 0 \cdot 9 + \\ &+ u_{32}d(x_3, x_2) + u_{35}d(x_3, x_5) + 1 \cdot 1 + 0 \cdot 10 + \\ &+ u_{42}d(x_4, x_2) + u_{45}d(x_4, x_5) + 0 \cdot 10 + 1 \cdot 1 + \\ &+ u_{52}d(x_5, x_2) + u_{55}d(x_5, x_5) + 0 \cdot 9 + 1 \cdot 0 + \\ &+ u_{62}d(x_6, x_2) + u_{65}d(x_6, x_5) + 0 \cdot 10 + 1 \cdot 1 \end{aligned} = 4$$

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The PAM algorithm (Example)

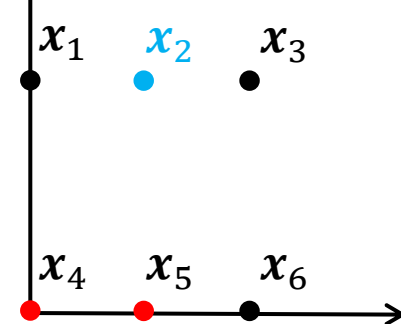
**Data set:**  $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ , with

$x_1 = [0,3]^T$ ,  $x_2 = [1,3]^T$ ,  $x_3 = [2,3]^T$ ,  $x_4 = [0,0]^T$ ,  $x_5 = [1,0]^T$ ,  $x_6 = [2,0]^T$

**Sets of medoids:**  $\Theta = \{x_4, x_5\}$ ,  $\Theta_{42} = \{x_2, x_5\}$

Computation of  $\Delta J_{42}$  as

$\Delta J_{42} = J(\Theta_{42}, U_{42}) - J(\Theta, U) = \sum_{h \in X - \Theta} C_{h42}$  (Sq. Eucl. dist. is used):



	Dist. from Closest repr. in $\Theta = \{x_4, x_5\}$	Dist. from Next closest repr. in $\Theta = \{x_4, x_5\}$	Dist. from closest repr. in $\Theta_{42} = \{x_2, x_5\}$	$C_{h42}$
$x_1$	9 ( $x_4$ )	10 ( $x_5$ )	1 ( $x_2$ )	$1 - 9 = -8$
$x_2$	9 ( $x_5$ )	10 ( $x_4$ )	0 ( $x_2$ )	$0 - 9 = -9$
$x_3$	10 ( $x_5$ )	13 ( $x_4$ )	1 ( $x_2$ )	$1 - 10 = -9$
$x_4$	0 ( $x_4$ )	1 ( $x_5$ )	1 ( $x_5$ )	$1 - 0 = 1$
$x_5$	0 ( $x_5$ )	1 ( $x_4$ )	0 ( $x_5$ )	$0 - 0 = 0$
$x_6$	1 ( $x_5$ )	2 ( $x_4$ )	1 ( $x_5$ )	$2 - 1 = 1$
$\Delta J_{42}$				<b><math>-25</math></b>

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The CLARA algorithm

- It is more suitable for large data sets.
- **The strategy:**
  - **Draw** randomly a **sample  $X'$**  of size  $N'$  from the entire data set.
  - **Run** the **PAM** algorithm to **determine  $\theta'$**  that best represents  $X'$ .
  - Use  $\theta'$  in the place of  $\theta$  to represent the entire data set  $X$ .
- **The rationale:**
  - Assuming that  $X'$  has been selected in a way **representative** of the **statistical distribution** of the **data points** in  $X$ ,  $\theta'$  is expected to be a good approximation of  $\theta$ , which would have been produced if PAM were run on the entire  $X$ .
- **The algorithm:**
  - Draw  $s$  sample subsets of size  $N'$  from  $X$ , denoted by  $X'_1, \dots, X'_s$  (typically  $s = 5, N' = 40 + 2m$ ).
  - Run PAM on each one of them and identify  $\theta'_1, \dots, \theta'_s$ .
  - Choose the set  $\theta'_j$  that minimizes

$$J(\theta, U) = \sum_{i \in I_{X-\theta'}} \sum_{j \in I_{\theta'}} u_{ij} d(\mathbf{x}_i, \mathbf{x}_j)$$

based on the entire data set  $X$ .



# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The CLARANS algorithm

- It is more **suitable** for **large data sets**.
- It follows the philosophy of PAM with the difference that only **a randomly selected fraction**  $q (< m(N - m))$  of the **neighbors** of the current medoid set is **considered**.
- It performs several runs ( $s$ ) starting from different initial choices for  $\Theta$ .

### The algorithm:

- For  $i = 1$  to  $s$ 
  - o **Initialize** randomly  $\Theta$ .
  - o **(A) Select** randomly  $q$  neighbors of  $\Theta$ .
  - o For  $j = 1$  to  $q$ 
    - \* **If** the present **neighbor of  $\Theta$**  is **better** than  $\Theta$  (in terms of  $J(\Theta, U)$ ) then
      - **Set  $\Theta$**  equal to **its neighbor**
      - Go to **(A)**
    - \* **End If**
  - o **End For**
  - o Set  $\Theta^i = \Theta$
- **End For**
- **Select** the **best  $\Theta^i$**  with respect to  $J(\Theta, U)$ .
- Based on  $\Theta^i$ , **assign** each  $x \in X - \Theta$  to the cluster whose representative is closest to  $x$

# CFO hard clustering algorithms

## Generalized Hard Algorithmic Scheme (GHAS)

### The CLARANS algorithm (cont.)

#### Remarks:

- **CLARANS depends** on  $q$  and  $s$ . Typically,  $s = 2$  and
$$q = \max(0.125m(N - m), 250)$$
- As  $q$  approaches  $m(N - m)$  CLARANS approaches PAM and the complexity increases.
- CLARANS can also be described in terms of graph theory concepts.
- **CLARANS unravels better quality clusters than CLARA.**
- In some cases, CLARA is significantly faster than CLARANS.
- **CLARANS retains its quadratic computational nature** and thus it is not appropriate for very large data sets.

# Probability and statistics: a brief review

**Random variable (RV):** It models the output of an experiment.

## RV types:

- Discrete
- continuous

## Discrete random variables:

- A **discrete RV**  $x$  can take any value  $x$  from a **finite** or **countably infinite** set  $X$ .
- $X$ : **sample space** or **state space**.
- **Event:** Any **subset** of  $X$ .
- **Elementary or simple event:** A **single element subset** of  $X$ .
- **Example:** Consider the die roll experiment.  $X = \{1, 2, 3, 4, 5, 6\}$
- Events: “Odd number”, “number  $> 3$ ”, “2”, “5”

Elementary events

# Probability and statistics: a brief review

## Discrete random variables (cont.):

- **Notation:** **Probability** of the **event**  $x=x \in X$ :  $P(x=x) \equiv P(x)$
- $P(\cdot)$ : A function called **probability mass function (pmf)** satisfying
  - ✓  $P(x) \geq 0, \forall x \in X$
  - ✓  $\sum_{x \in X} P(x) = 1$

# Probability and statistics: a brief review

## Discrete random variables (cont.):

The case of more than one random variables: **Definitions**

Discrete RV	$x$	$y$
Sample space	$X=\{x_1, \dots, x_{nx}\}$	$Y=\{y_1, \dots, y_{ny}\}$

**Joint probability:**  $P(x_i, y_j) \equiv P(x=x_i \text{ AND } y=y_j)$

- It corresponds to the case where  $x$  takes the value  $x_i$  **AND**  $y$  takes the value  $y_j$ , **simultaneously**.

**Marginal probabilities:**  $P(x_i) \equiv P(x=x_i)$ ,  $P(y_j) = P(y=y_j)$

- This terminology is used only when more than one rvs are involved.

**Conditional probability:**  $P(x_i | y_j) \equiv P(x=x_i | y=y_j) = P(x_i, y_j) / P(y_j)$

- It corresponds to the case where  $x$  takes the value  $x_i$  **given that**  $y$  takes the value  $y_j$ .

# Probability and statistics: a brief review

## Discrete random variables (cont.):

The case of more than one variables: *Properties*

Discrete RV	$x$	$y$
Sample space	$X = \{x_1, \dots, x_{nx}\}$	$Y = \{y_1, \dots, y_{ny}\}$

**Sum rule:**  $P(x) = \sum_{y \in Y} P(x, y), \quad \forall x \in X$

**Product rule:**  $P(x, y) = P(x | y)P(y)$

**Statistical independence:**  $P(x, y) = P(x)P(y)$

A consequence:  $P(x | y) = P(x) \quad P(y | x) = P(y)$

**Bayes rule:**  $P(y | x) = \frac{P(x | y)P(y)}{P(x)}$

It plays a **key role** in **ML**.

or

$$P(y | x) = \frac{P(x | y)P(y)}{\sum_{y \in Y} P(x | y)P(y)}$$

# Probability and statistics: a brief review

## Continuous random variables:

- A **continuous RV**  $x$  can take any value  $x \in R$ .

- **Sample space** or **state space**:  $R$

- **Events**:  $\{x \leq x\}$ ,  $\{x_1 < x \leq x_2\}$ ,  $\{x \geq x\}$

- **Cumulative distribution function (cdf)**:  $F_x(x) = P(x \leq x)$

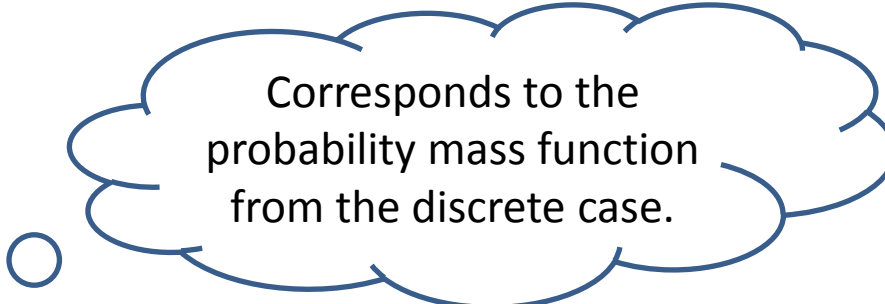
- It is  $F_x(\infty) = P(x < \infty) = 1$

- **Probability of events** in terms of **cdf**:

- $P(x \leq x) = F_x(x)$

- $P(x_1 < x \leq x_2) = P(x \leq x_2) - P(x \leq x_1) = F_x(x_2) - F_x(x_1)$

- $P(x \geq x) = P(x \leq \infty) - P(x \leq x) = 1 - P(x \leq x) = 1 - F_x(x)$



Corresponds to the probability mass function from the discrete case.



It assigns "mass" to events.

# Probability and statistics: a brief review

## Continuous random variables (cont.):

• **Assumption:**  $F_x(x)$  is *continuous* and *differentiable*.

• **Probability density function (pdf):**

$$p_x(x) = \frac{dF_x(x)}{dx}$$

It assigns “mass” to values.

• **cdf in terms of pdf:**

$$F_x(x) = \int_{-\infty}^x p_x(z) dz$$

• **Probability of events in terms of pdf:**

$$\blacktriangleright P(x \leq x) = F_x(x) = \int_{-\infty}^x p_x(z) dz$$

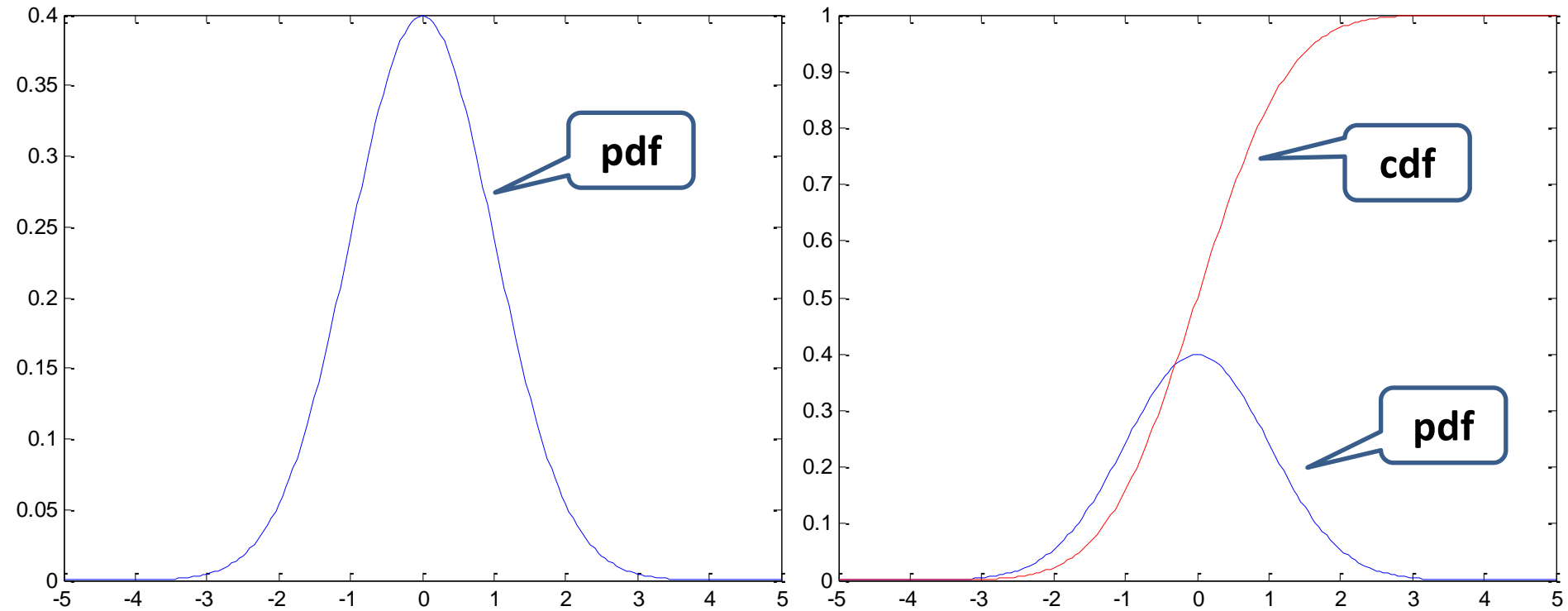
$$\blacktriangleright P(x_1 < x \leq x_2) = P(x \leq x_2) - P(x \leq x_1) = F_x(x_2) - F_x(x_1) = \int_{x_1}^{x_2} p_x(x) dx$$

$$\blacktriangleright P(x \geq x) = P(x \leq \infty) - P(x \leq x) = 1 - P(x \leq x) = 1 - F_x(x) = \int_{-\infty}^x p_x(z) dz$$



# Probability and statistics: a brief review

## Continuous random variables (cont.):



# Probability and statistics: a brief review

## Continuous random variables (cont.):

• Since  $P(-\infty < x < +\infty) = 1$  it is:  $\int_{-\infty}^{+\infty} p_x(x) dx = 1$

• It is  $P(x < x \leq x + \Delta x) = \int_x^{x+\Delta x} p_x(z) dz \approx p_x(x) \Delta x$

As  $\Delta x \rightarrow 0$ ,  $P(x < x < x + \Delta x) = P(x = x) = 0$ .

The **probability** of a **continuous rv** to take a **single value** is **zero**.

## The case of more than one variables:

Continuous RV	$x$	$y$
Sample space	$R$	$R$

**NOTE:** All rules stated for the **probability mass function** in the **discrete case** are stated for the **pdf** in the **continuous case**.

### Product rule

$$p(x, y) = p(x | y) p(y)$$

We drop the name of rv from the subscript of  $p$ .

### Sum rule

$$p(x) = \int_{-\infty}^{+\infty} p(x, y) dy$$

# Probability and statistics: a brief review

## Useful quantities related to (continuous) rvs:

For **discrete** rv's, the integrals become summations.

• Mean (expected) value of a rv  $x$ :  $E[x] = \int_{-\infty}^{+\infty} xp(x)dx$

• Variance of a rv  $x$ :  $\sigma_x^2 = \int_{-\infty}^{+\infty} (x - E[x])^2 p(x)dx = E[(x - E(x))^2]$

• Mean (expected) value of a function of an rv  $x$ :  $E[f(x)] = \int_{-\infty}^{+\infty} f(x)p(x)dx$

• Mean of a function of two rv's  $x, y$ :  $E_{x,y}[f(x, y)] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f(x, y)p(x, y)dxdy$

• Conditional mean of an rv  $y$  given  $x = x$ :  $E[y | x] = \int_{-\infty}^{+\infty} yp(y | x)dy$

• It is  $E_{x,y}[f(x, y)] = E_x[E_{y|x}[f(x, y)]]$

• Covariance between two rvs  $x$  and  $y$ :  $\text{cov}(x, y) = E[(x - E[x])(y - E[y])]$

• Correlation between two rv's  $x$  and  $y$ :  $r_{xy} \equiv E(xy) = \text{cov}(x, y) + E[x]E[y]$

• Correlation coefficient  $r_{xy} = \frac{E[x - E[x]](y - E[y])}{\sigma_x \sigma_y}$

# Probability and statistics: a brief review

## Random vectors

• A collection of rvs:  $\mathbf{x} = [x_1, x_2, \dots, x_l]^T$

• Probability density function (pdf) of  $\mathbf{x}$ : The joint pdf of  $x_1, x_2, \dots, x_l$ .  
 $p(\mathbf{x}) = p(x_1, x_2, \dots, x_l)$

• Covariance matrix of  $\mathbf{x}$ :

$$\text{cov}(\mathbf{x}) = E[(\mathbf{x} - E[\mathbf{x}])(\mathbf{x} - E[\mathbf{x}])^T] = \begin{bmatrix} \text{cov}(x_1, x_1) & \cdots & \text{cov}(x_1, x_l) \\ \vdots & \ddots & \vdots \\ \text{cov}(x_l, x_1) & \cdots & \text{cov}(x_l, x_l) \end{bmatrix}$$

• Correlation matrix of  $\mathbf{x}$ :  $R_{\mathbf{x}} = E[\mathbf{x}\mathbf{x}^T] = \begin{bmatrix} E(x_1 x_1) & \cdots & E(x_1 x_l) \\ \vdots & \ddots & \vdots \\ E(x_l x_1) & \cdots & E(x_l x_l) \end{bmatrix}$

• It is  $R_{\mathbf{x}} \equiv E[\mathbf{x}\mathbf{x}^T] = \text{cov}(\mathbf{x}) + E[\mathbf{x}]E[\mathbf{x}^T]$

**Exercise:** Prove this identity

# Probability and statistics: a brief review

## Random vectors (cont.)

- Remark: Both  $R_{\mathbf{x}}$  and  $\text{cov}(\mathbf{x})$  are **symmetric** and **positive definite**  $l \times l$  matrices.

**Exercise:** Prove these statements

A square matrix  
 $A$  is **symmetric**  
iff  $A^T = A$ .

A square matrix  
 $A$  is **positive  
definite** iff  
 $\mathbf{z}^T A \mathbf{z} > 0, \forall \mathbf{z} \in \mathbb{R}^l$ .

# Probability and statistics: a brief review

- **One dim. normal (Gaussian) distribution**  $x \sim N(\mu, \sigma^2)$  or  $N(x | \mu, \sigma^2)$  :

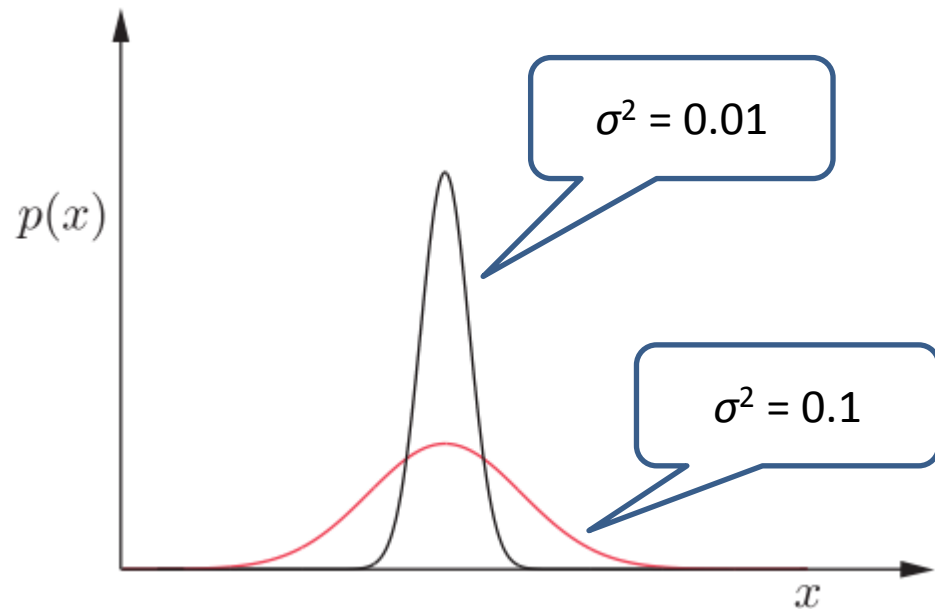
- Sample space:  $\mathcal{R}$

- It is

- $$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

- $E[x] = \mu$

- $\sigma_x^2 = \sigma^2.$



# Probability and statistics: a brief review

- Multi dim. normal (Gaussian) distribution  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  or  $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ :

- $l$ -dim. case

- It is

$$\triangleright p(\mathbf{x}) = \frac{1}{(2\pi)^{l/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}{2}\right)$$

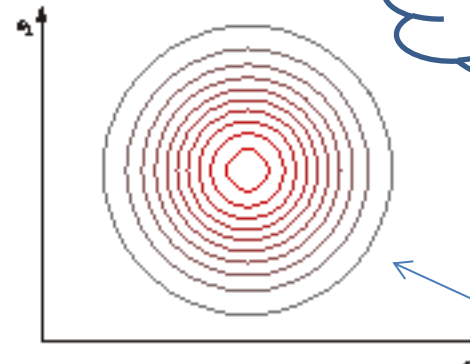
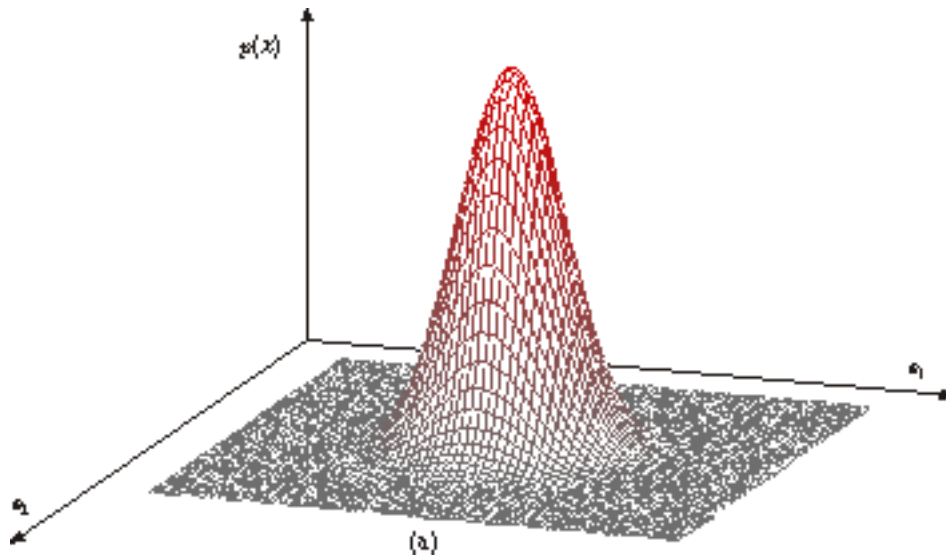
$$\triangleright E[\mathbf{x}] = \boldsymbol{\mu}$$

$$\triangleright \text{cov}(\mathbf{x}) = \boldsymbol{\Sigma}.$$

(\*) For the 2-d case  $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$

# Probability and statistics: a brief review

- Multi dim. normal (Gaussian) distribution  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  or  $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ :

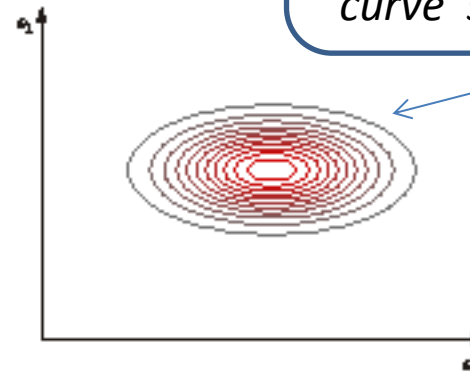
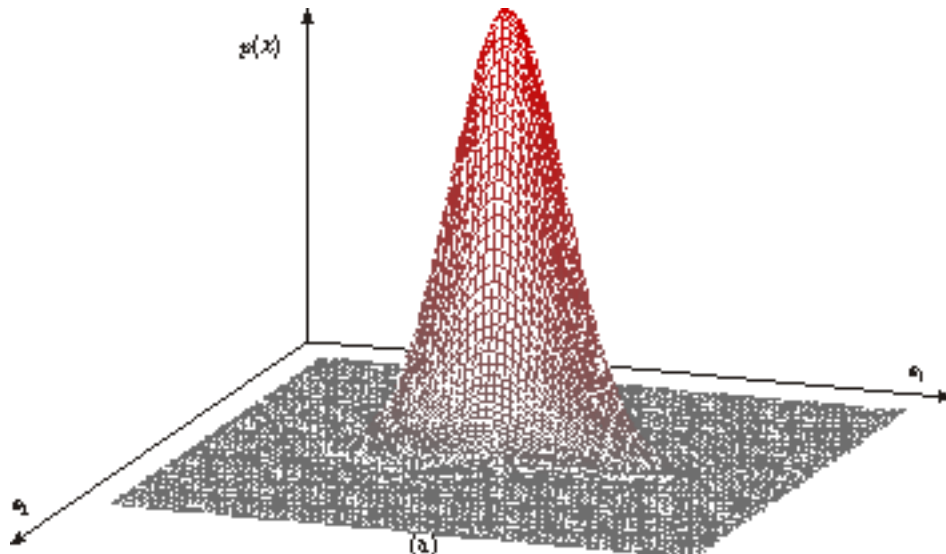


$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$$

$\boldsymbol{\Sigma}$ : diagonal with  $\sigma_1^2 = \sigma_2^2$

**Isovalued curves:**

- $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \text{const.}$
- All points on each isovalue curve share the value  $p(\mathbf{x})$ .

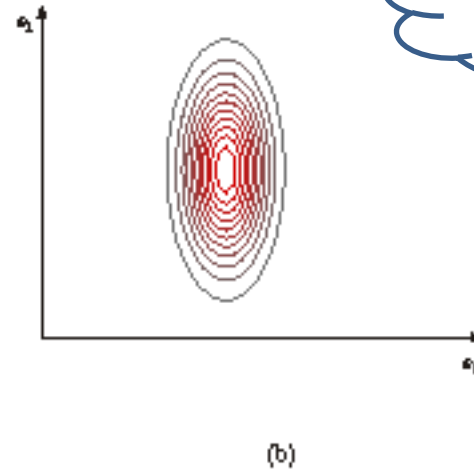
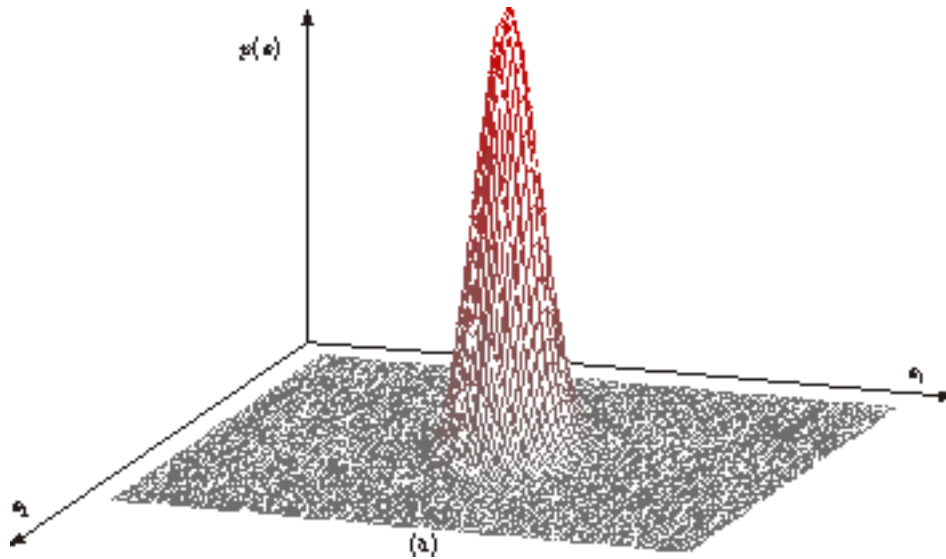


$\boldsymbol{\Sigma}$ : diagonal with  $\sigma_1^2 \gg \sigma_2^2$



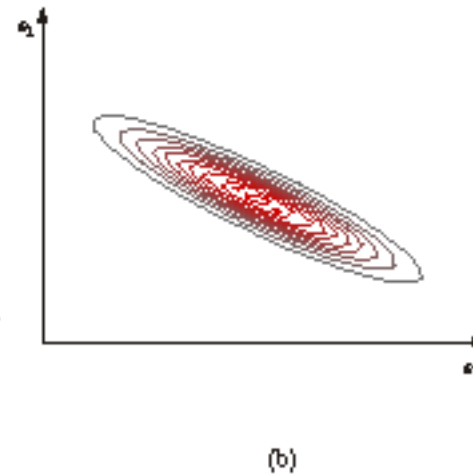
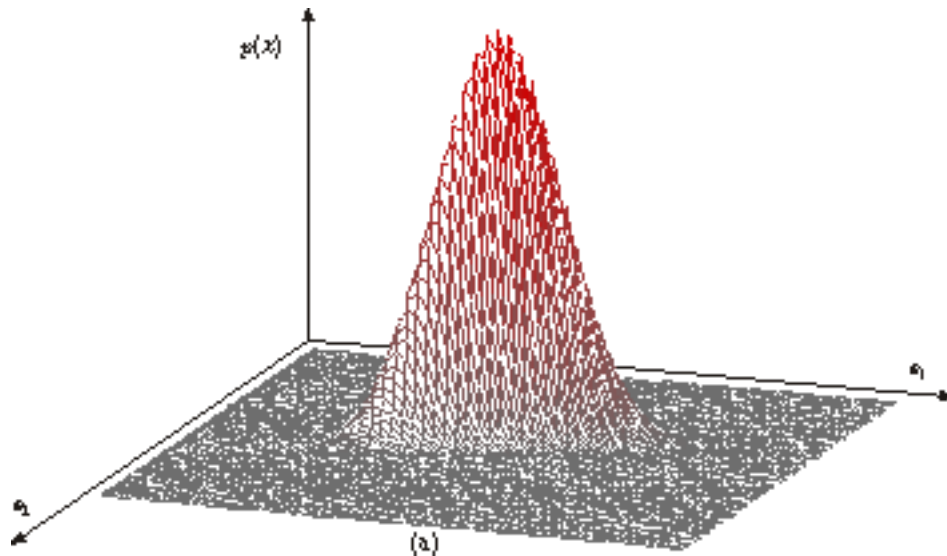
# Probability and statistics: a brief review

- Multi dim. normal (Gaussian) distribution  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  or  $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ :



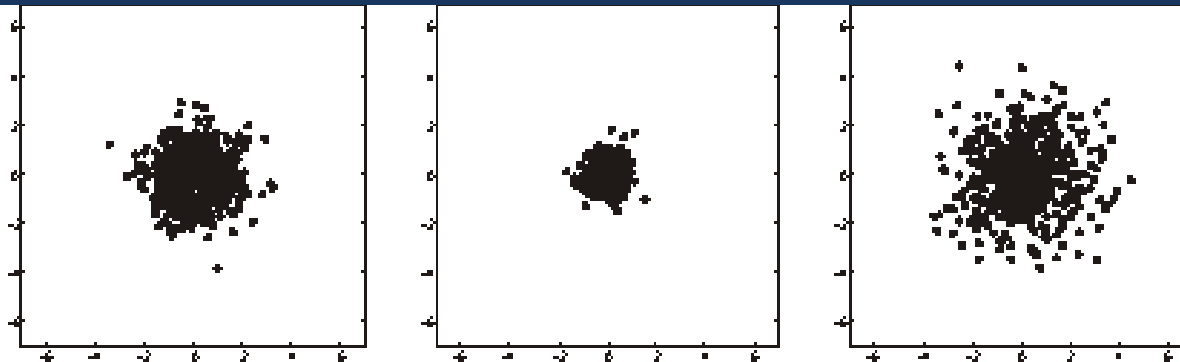
$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$$

$\boldsymbol{\Sigma}$ : diagonal  
with  $\sigma_1^2 \ll \sigma_2^2$



$\boldsymbol{\Sigma}$ : non diagonal

# Probability and statistics: a brief review



(a)

(b)

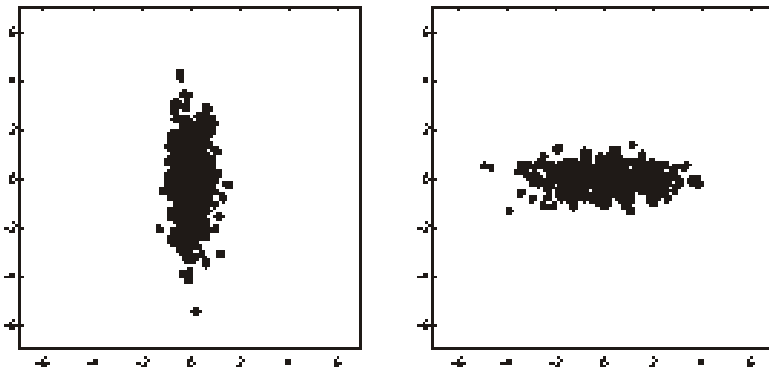
(c)

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$$

(a)  $\sigma_1^2 = \sigma_2^2 = 1, \sigma_{12} = 0$

(b)  $\sigma_1^2 = \sigma_2^2 = 0.2, \sigma_{12} = 0$

(c)  $\sigma_1^2 = \sigma_2^2 = 2, \sigma_{12} = 0$

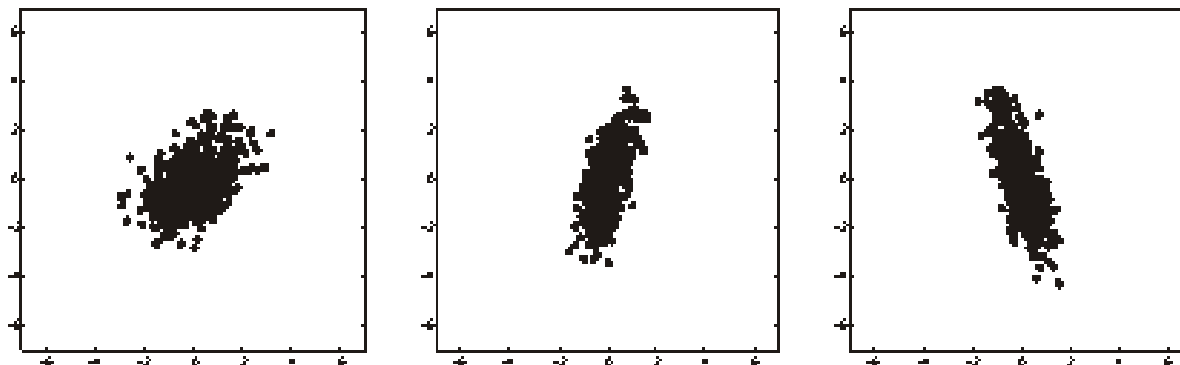


(d)

(e)

(d)  $\sigma_1^2 = 0.2, \sigma_2^2 = 2, \sigma_{12} = 0$

(e)  $\sigma_1^2 = 2, \sigma_2^2 = 0.2, \sigma_{12} = 0$



(f)

(g)

(h)

(f)  $\sigma_1^2 = \sigma_2^2 = 1, \sigma_{12} = 0.5$

(g)  $\sigma_1^2 = 0.3, \sigma_2^2 = 2, \sigma_{12} = 0.5$

(h)  $\sigma_1^2 = 0.3, \sigma_2^2 = 2, \sigma_{12} = -0.5$

# Probability and statistics: a brief review

## Continuous RV distributions (cont.)

### • Multi dim. normal (Gaussian) distribution $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ or $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ :

From 1-dim.  $\rightarrow$  2-dim. case.

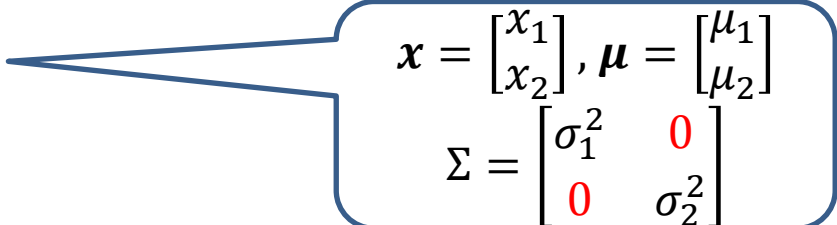
▪ **1-dim. case:**  $p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = \frac{1}{(2\pi)^{1/2}\sigma} \exp\left(-\frac{(x-\mu)\sigma^{-2}(x-\mu)}{2}\right)$

▪ A **first extension** to the 2-dim. case (**independent** rv's):

▪  $p(x_1, x_2) = p_1(x_1) \cdot p_2(x_2) =$

▪  $\frac{1}{(2\pi)^{1/2} \cdot \sigma_1} \exp\left(-\frac{(x_1-\mu_1)\sigma_1^{-2}(x_1-\mu_1)}{2}\right) \cdot \frac{1}{(2\pi)^{1/2} \cdot \sigma_2} \exp\left(-\frac{(x_2-\mu_2)\sigma_2^{-2}(x_2-\mu_2)}{2}\right) =$

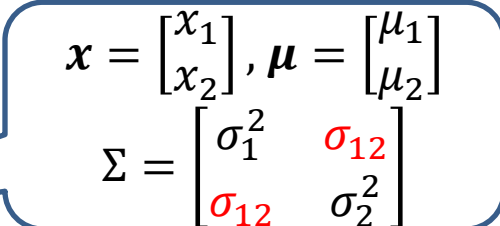
$\frac{1}{(2\pi)|\Sigma|^{1/2}} \exp\left(-\frac{(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}{2}\right)$



$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$   
 $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$

▪ The **final extension** to the 2-dim. case (**dependent** rv's):

▪  $p(x_1, x_2) = \frac{1}{(2\pi)|\Sigma|^{1/2}} \exp\left(-\frac{(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}{2}\right)$



$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$   
 $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$

# Probability and statistics: a brief review

## • Multi dim. normal (Gaussian) distribution $x \sim N(\mu, \Sigma)$ or $N(x | \mu, \Sigma)$ :

### ▪ Properties

1. If the covariance matrix  $\Sigma$  is **diagonal**, then, the **rv's  $x_1, \dots, x_l$**  comprising  **$x$**  are **statistically independent**. It is

$$p(\mathbf{x}) = \prod_{i=1}^l p_i(x_i) = \prod_{i=1}^l \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right)$$

### 2. Central limit theorem:

Let:

- $x_1, \dots, x_r$  independent rvs following different distributions
- $\mu_i, \sigma_i^2$  mean and variance of  $x_i$ .
- Define  $x = x_1 + \dots + x_r$ ,  $\mu = \mu_1 + \dots + \mu_r$ ,  $\sigma^2 = \sigma_1^2 + \dots + \sigma_r^2$ .
- Define  $z = (x - \mu)/\sigma$ .

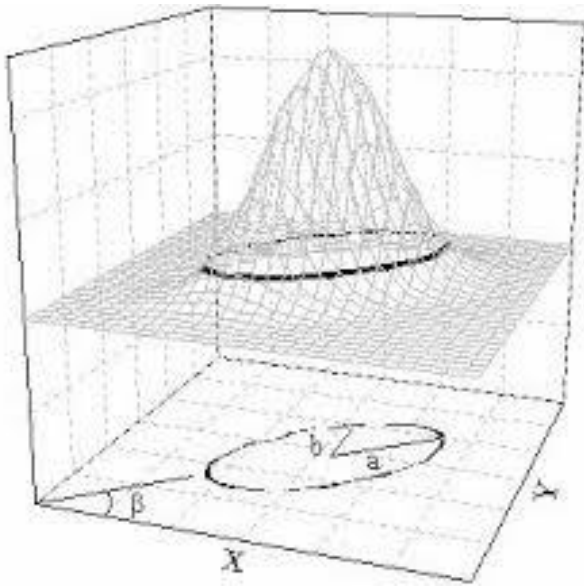
Then

- $p(z) \rightarrow N(z|0,1)$ , as  $r \rightarrow \infty$

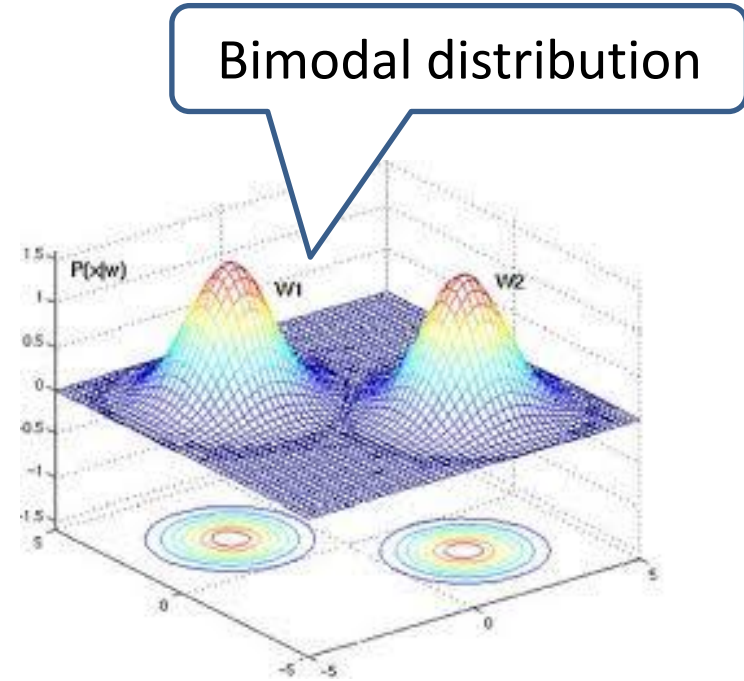
# Probability and statistics: a brief review

## Continuous RV distributions (cont.)

### ▪ Other examples of multi-dimensional pdfs



Two-dim. pdfs



# Probability and statistics: a brief review

## Likelihood function

- Let  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  a set of independent data vectors
- Let  $p_{\theta}(\cdot)$  be a pdf belonging to a **known parametric set of pdf functions** of parameter vector  $\theta$ .
- $p(\mathbf{x}) = p_{\theta}(\mathbf{x}) \equiv p(\mathbf{x}; \theta)$ .

### Examples:

- If  $p_{\theta}(\mathbf{x})$  is **normal** distribution **parameterized** on the mean vector  $\mu$ ,  $\theta$  will simply be  $\mu$ .
- If  $p_{\theta}(\mathbf{x})$  is **normal** distribution **parameterized** on both the mean vector  $\mu$  and the cov. matrix  $\Sigma$ ,  $\theta$  will contain the coordinates of both  $\mu$  and  $\Sigma$ .

**Likelihood function** of  $\theta$  wrt  $X$ :  $p(X; \theta) = p(\mathbf{x}_1, \dots, \mathbf{x}_N; \theta) = \prod_{i=1}^N p(\mathbf{x}_i; \theta)$

**Log-likelihood function** of  $\theta$  wrt  $X$ :

$$L(\theta) = \ln p(X; \theta) = \ln p(\mathbf{x}_1, \dots, \mathbf{x}_N; \theta) = \sum_{i=1}^N \ln p(\mathbf{x}_i; \theta)$$

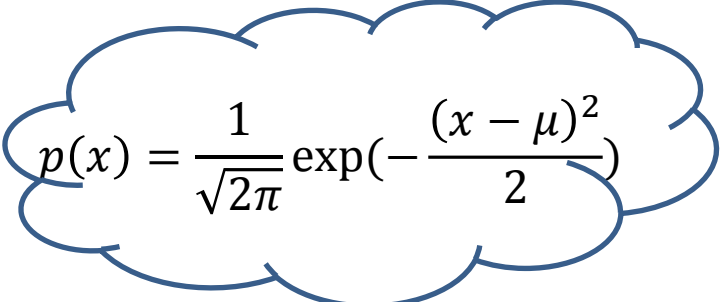
# Probability and statistics: a brief review

## Likelihood function

### Example:

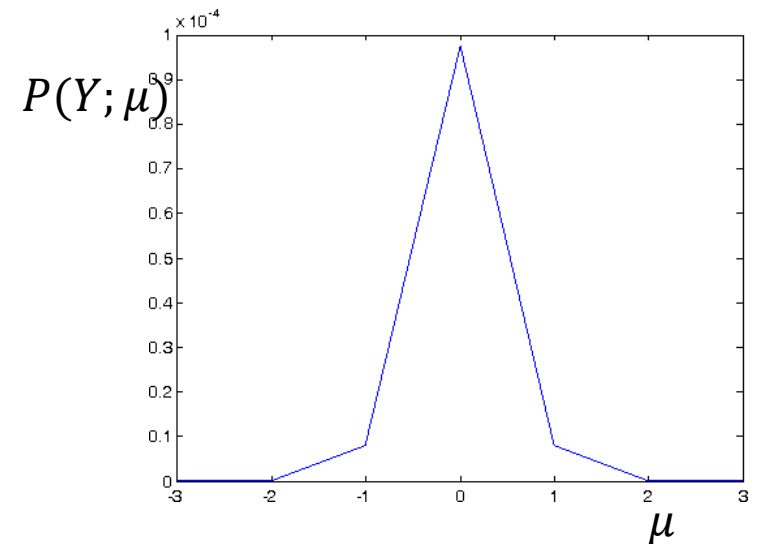
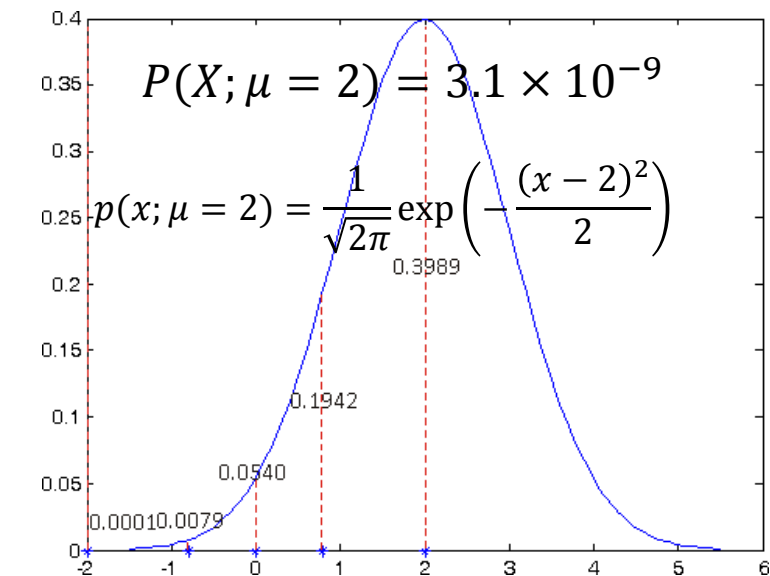
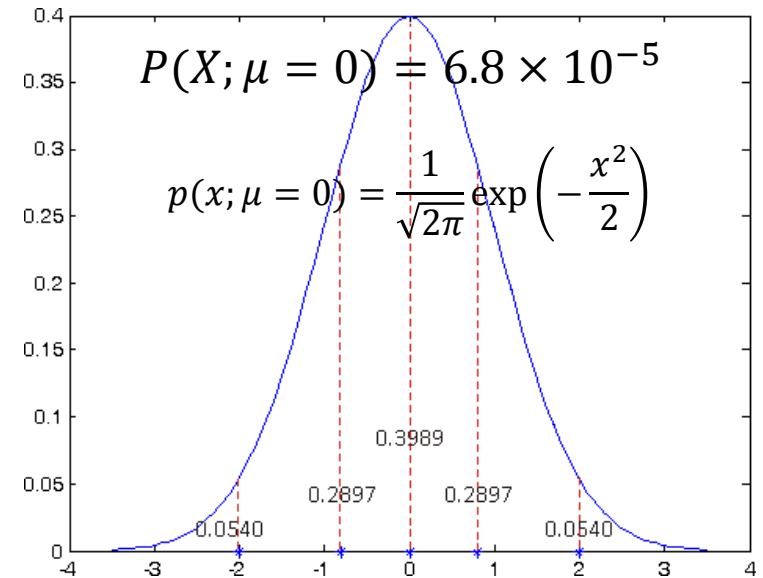
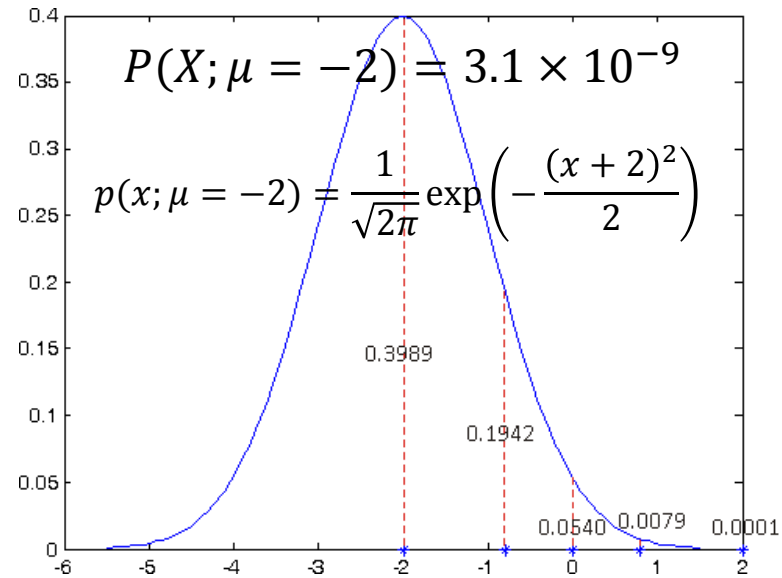
- $X = \{-2, -1, 0, 1, 2\}$
- Consider the **parametric set** of **normal distributions** of **unit variance**, parameterized on  $\mu$ .
- The likelihood of  $\mu$  wrt  $X$  is

$$p(X; \mu) = p(-2, -1, 0, 1, 2; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(-2-\mu)^2}{2}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(-1-\mu)^2}{2}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(0-\mu)^2}{2}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(1-\mu)^2}{2}\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(2-\mu)^2}{2}\right)$$


$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2}\right)$$

# Probability and statistics: a brief review

## Likelihood function





# Probabilistic CFO clustering algorithms

## Maximum likelihood (ML) method:

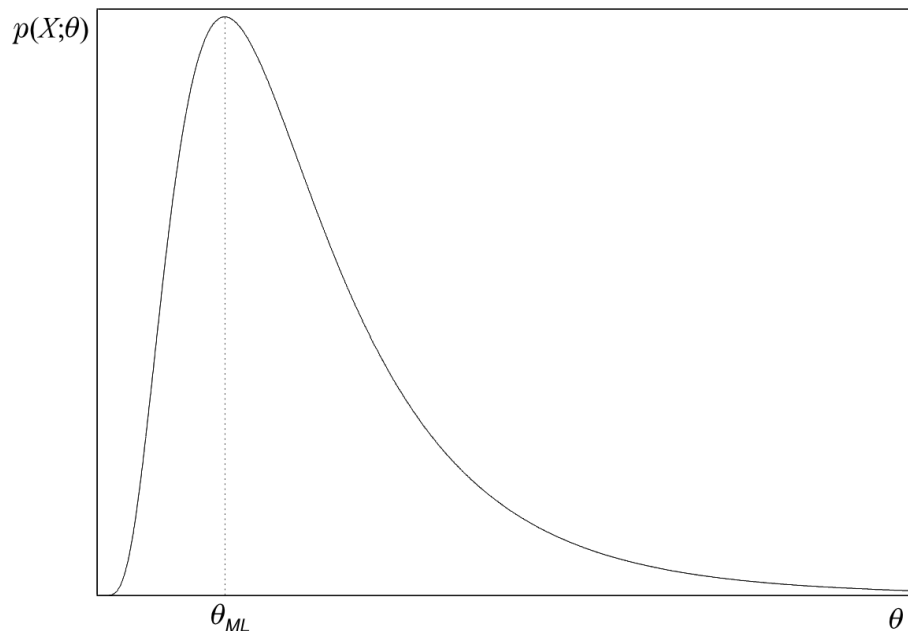
Given a set of independent data vectors  $Y = \{x_1, x_2, \dots, x_N\}$ ,

**estimate** the parameter vector  $\theta$  as the **maximum** of the **likelihood** ( $p(Y; \theta)$ ) or the **log-likelihood** ( $L(\theta)$ ) function.

$$\hat{\theta}_{ML} = \operatorname{argmax}_{\theta} p(Y; \theta)$$

→

$$\hat{\theta}_{ML}: \frac{\partial L(\theta)}{\partial \theta} = \sum_{k=1}^N \frac{1}{p(x_k; \theta)} \frac{\partial p(x_k; \theta)}{\partial \theta} = \mathbf{0}$$



Since  $\ln(\cdot)$  is an **increasing function**,  $p(Y; \theta)$  and  $L(\theta)$  share the **same maxima**.

# Probabilistic CFO clustering algorithms

## Maximum likelihood (ML) method:

### Assuming that

- the chosen model  $p(\mathbf{x}; \boldsymbol{\theta})$  is **correct** and
- there **exists** a true parameter  $\boldsymbol{\theta}_o$ ,

### the ML estimator

- (a) is **asymptotically unbiased**  $\lim_{N \rightarrow \infty} E[\hat{\boldsymbol{\theta}}_{ML}] = \boldsymbol{\theta}_o$
- (b) is **asymptotically consistent**  $\lim_{N \rightarrow \infty} Prob\{\|\hat{\boldsymbol{\theta}}_{ML} - \boldsymbol{\theta}_o\|\} = 0$
- (c) is **asymptotically efficient** (it achieves the **Cramer-Rao lower bound**)

The **pdf** of the ML estimator **approaches** the **normal distribution** with mean  $\boldsymbol{\theta}_o$ , as  $N \rightarrow \infty$ .

# Maximum likelihood method

## Example 1:

-Let  $Y$  be a set of  $N$  (independent from each other) data points,  $\mathbf{x}_i$ ,  $i = 1, \dots, N$ , generated by a normal distribution  $p(\mathbf{x}; \boldsymbol{\theta})$  of known covariance matrix and unknown mean.

-Determine the ML estimate of the mean  $\boldsymbol{\mu}$  of  $p(\mathbf{x}; \boldsymbol{\theta})$ , based on  $Y$ .

## Solution:

-The unknown parameter vector in this case is the mean vector  $\boldsymbol{\mu}$ , i.e.  $\boldsymbol{\theta} \equiv \boldsymbol{\mu}$ .

-It is

$$p(\mathbf{x}; \boldsymbol{\theta}) \equiv p(\mathbf{x}; \boldsymbol{\mu}) = \frac{1}{(2\pi)^{l/2} |\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \Rightarrow$$

$$\ln p(\mathbf{x}; \boldsymbol{\mu}) = \ln \frac{1}{(2\pi)^{l/2} |\Sigma|^{1/2}} - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) = C - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

Then

$$L(\boldsymbol{\mu}) = \sum_{i=1}^N \ln p(\mathbf{x}_i; \boldsymbol{\mu}) = NC - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_i - \boldsymbol{\mu})$$

# Maximum likelihood method

## Example 1 (cont.):

Setting the **gradient** of  $L(\boldsymbol{\mu})$  wrt  $\boldsymbol{\mu}$  equal to  $\mathbf{0}$  we have

$$\frac{\partial L(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = \frac{\partial}{\partial \boldsymbol{\mu}} \left( NC - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) = \mathbf{0} \Leftrightarrow$$

$$\sum_{i=1}^N \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = \mathbf{0} \Leftrightarrow \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) = \mathbf{0} \Leftrightarrow \sum_{i=1}^N \mathbf{x}_i - N\boldsymbol{\mu} = \mathbf{0}$$

$$\boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

**Remark:** The **ML estimate** for the **covariance matrix** is

$$\boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$$

# Probabilistic CFO clustering algorithms

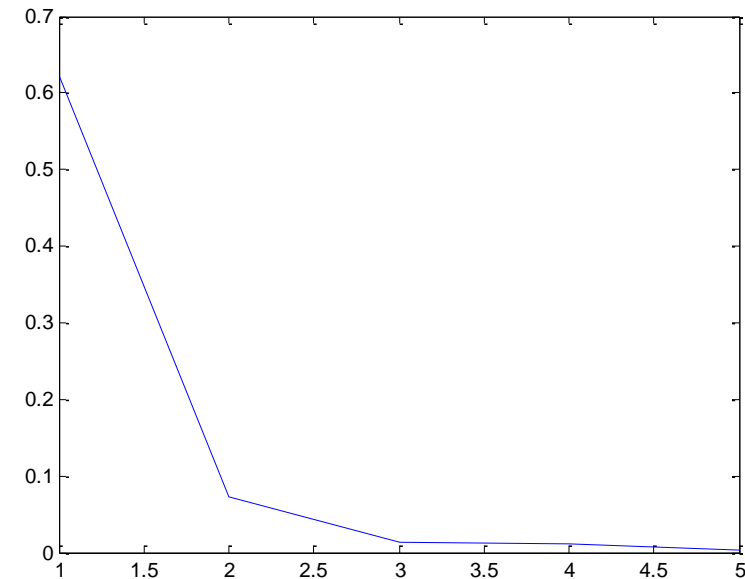
## Maximum likelihood (ML) method:

### Example:

- Let  $N = 10$  points resulting from the 1-dim. zero mean, unit variance normal distribution,  $N(0,1)$ .
- Let us pretend that we have at our disposal the following:
  - The knowledge that the distribution that generated the data is a normal one with variance equal to 1 and unknown mean vector  $\mu$ .
  - the 10 points.
- The ML estimate of the mean vector,  $\hat{\mu}$ , of the distribution is 0.6210 (the true value is 0).

The more the available data points, the more accurate the estimates for the mean vector.

$N$	Estimation error
10	0.6210 - 0
100	0.0727 - 0
1000	0.0138 - 0
10000	0.0118 - 0
100000	0.0034 - 0

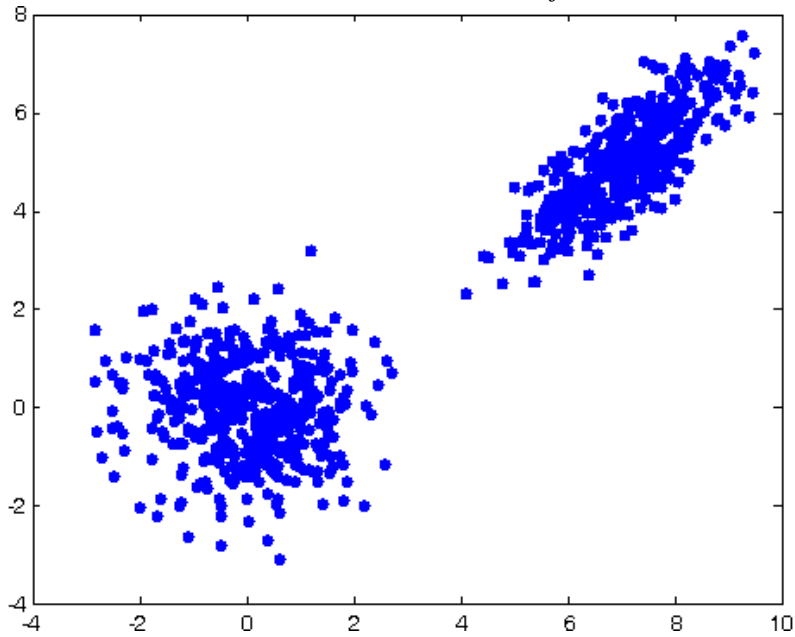


# Probabilistic CFO clustering algorithms

## Mixture models - The **Expectation – Maximization (EM)** algorithm

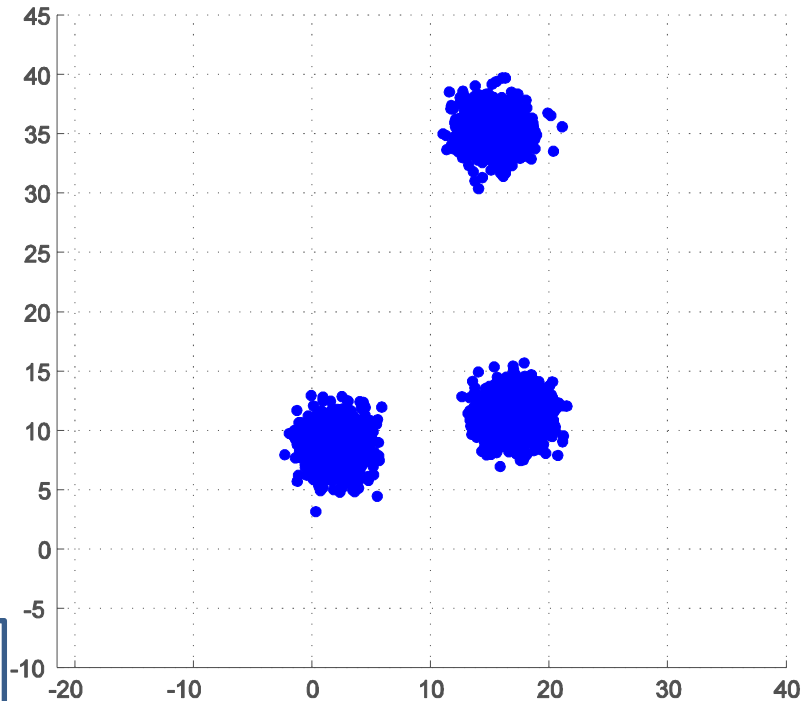
**Mixture model:** A **weighted sum** of **known parametric form pdfs**.

$$p(\mathbf{x}) = \sum_{j=1}^m P_j p(\mathbf{x} | j), \quad \sum_{j=1}^m P_j = 1, \quad \int_{-\infty}^{+\infty} p(\mathbf{x} | j) = 1$$



- Assume that  $p(\mathbf{x})$  models the distribution of the data in  $X$  (each pdf models a cluster).
- The **aim** is to **move** each pdf so that to **“cover”** the area in the data space where the vectors of each cluster lie (**mixture decomposition**).

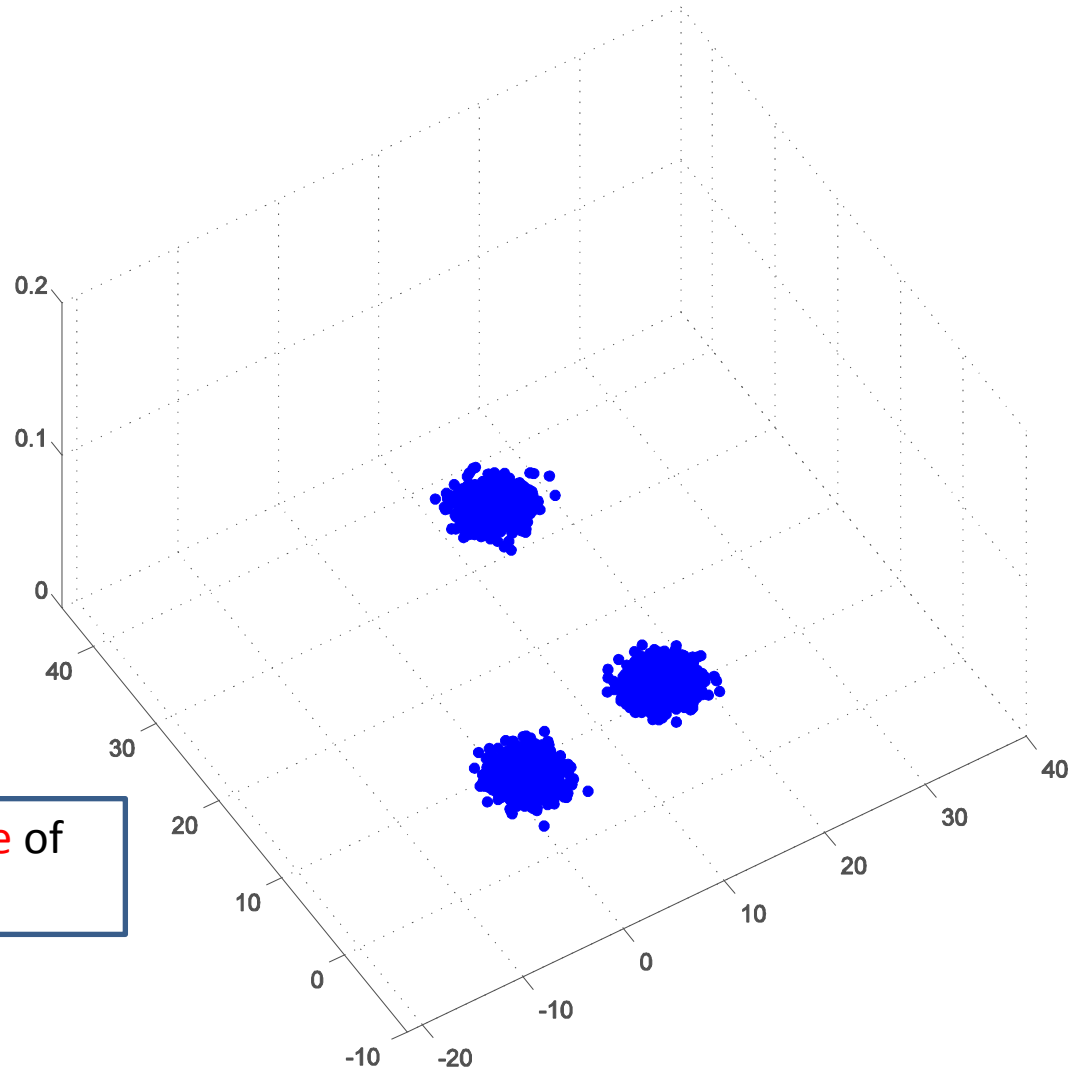
# Probabilistic CFO clustering algorithms



**Prerequisite:** Knowledge of the number of clusters.

- **Adopt** a **parametric mixture of distributions**, each one corresponding to a cluster (e.g., mixture of Gaussians), initialized randomly.
- **Move iteratively** the **distributions** each one above a **cluster**, **optimizing** a **criterion**.

# Probabilistic CFO clustering algorithms

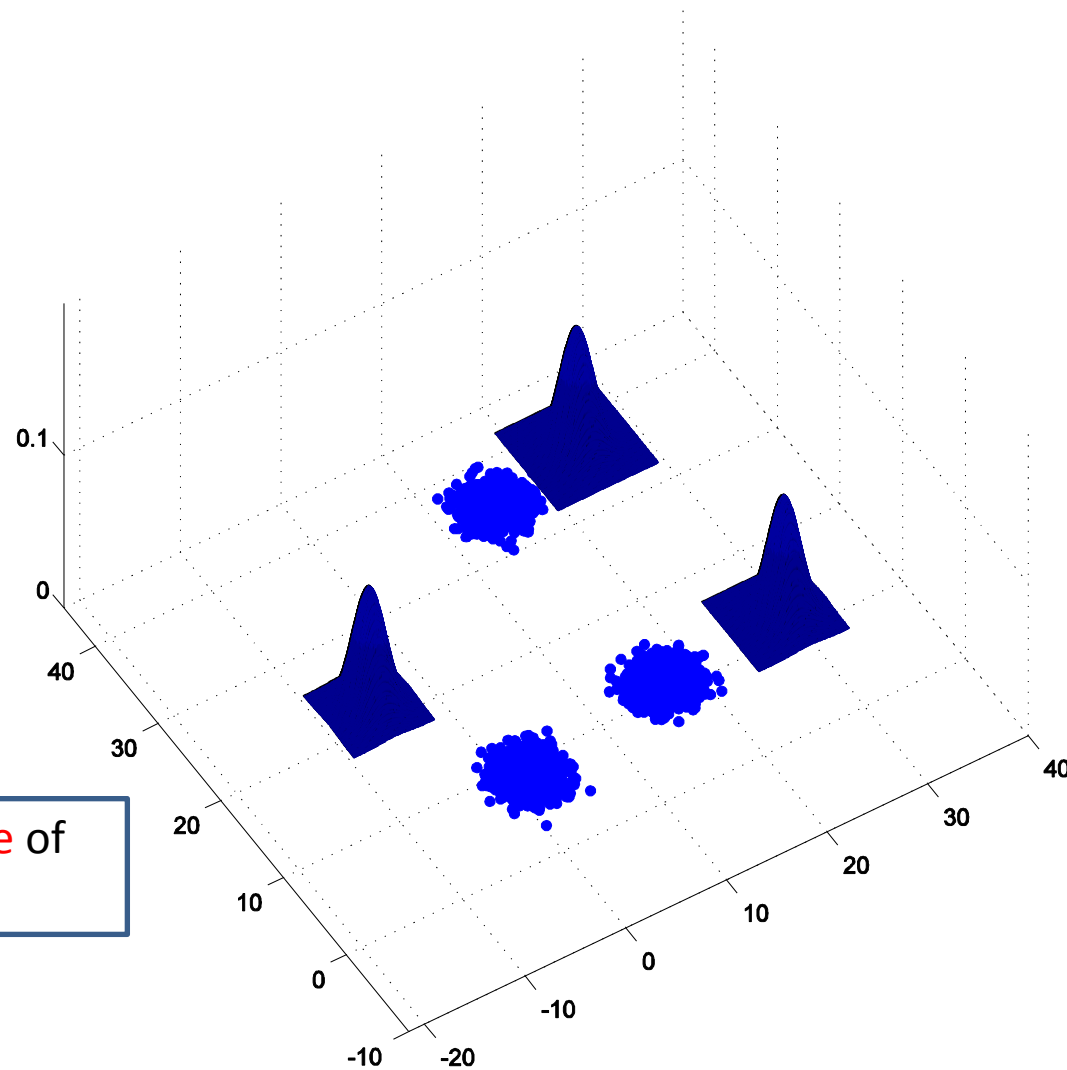


**Prerequisite:** Knowledge of the number of clusters.

- **Adopt** a **parametric mixture of distributions**, each one corresponding to a cluster (e.g., mixture of Gaussians), initialized randomly.
- **Move iteratively** the **distributions** each one above a **cluster**, **optimizing** a **criterion**.



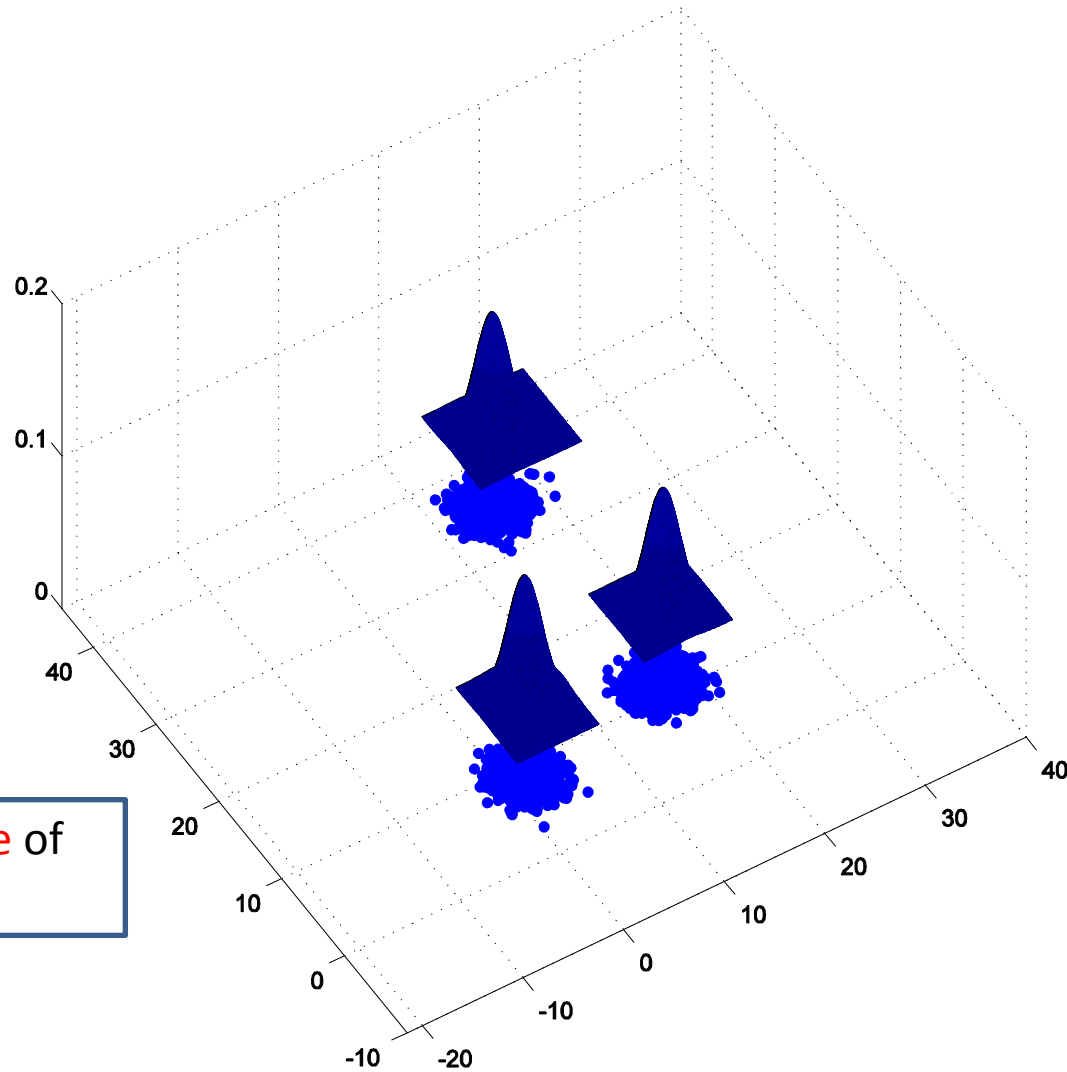
# Probabilistic CFO clustering algorithms



Prerequisite: Knowledge of the number of clusters.

- **Adopt** a **parametric mixture of distributions**, each one corresponding to a cluster (e.g., mixture of Gaussians), initialized randomly.
- **Move iteratively** the **distributions** each one above a **cluster**, **optimizing** a **criterion**.

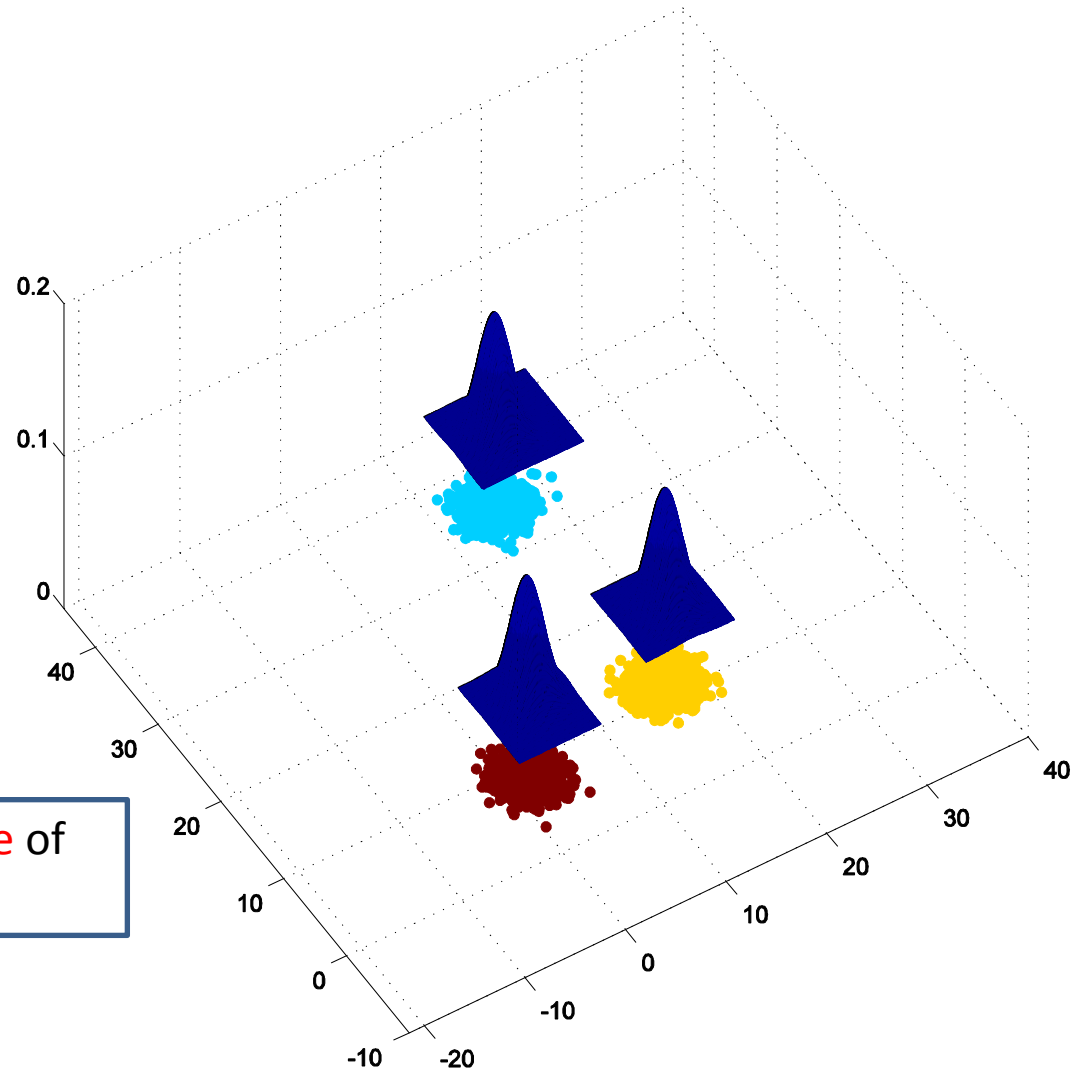
# Probabilistic CFO clustering algorithms



**Prerequisite:** Knowledge of the number of clusters.

- **Adopt** a **parametric mixture of distributions**, each one corresponding to a cluster (e.g., mixture of Gaussians), initialized randomly.
- **Move iteratively** the **distributions** each one above a **cluster**, **optimizing** a **criterion**.

# Probabilistic CFO clustering algorithms



**Prerequisite:** Knowledge of the number of clusters.

- **Adopt** a **parametric mixture of distributions**, each one corresponding to a cluster (e.g., mixture of Gaussians), initialized randomly.
- **Move iteratively** the **distributions** each one above a **cluster**, **optimizing** a **criterion**.

# Probabilistic CFO 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, with a **certain probability**.

Each **cluster** is **modeled** by a pdf  $p(\mathbf{x}|j)$ , parameterized by the vector  $\theta_j$ .

Let:

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

$P = \{P_1, P_2, \dots, P_m\}$ , the set of **a priori probabilities** of the **clusters**.

$P(j|\mathbf{x}) \equiv P(j|\mathbf{x}; \theta_j)$  the **(a posteriori) probability** of cluster  $j$ , given  $\mathbf{x}$ .

$p(\mathbf{x}|j) \equiv p(\mathbf{x}|j; \theta_j)$  the **pdf** that models cluster  $j$ .

It is  $p(\mathbf{x}) = \sum_{j=1}^m p(\mathbf{x}, j) = \sum_{j=1}^m p(\mathbf{x}|j) P_j$

**Bayes rule** 
$$P(j|\mathbf{x}) = \frac{p(\mathbf{x}, j)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|j)P_j}{p(\mathbf{x})}$$

# Probabilistic CFO clustering algorithms

It is

- $\sum_{j=1}^m P(j|\mathbf{x}_i) = 1, i = 1, \dots, N$
- $\sum_{j=1}^m P_j = 1.$

$$\text{ML: } L(\boldsymbol{\theta}) = \sum_{i=1}^N \ln(p(\mathbf{x}_i; \boldsymbol{\theta}))$$

Define the **cost function**

$$\begin{aligned} \ln p(X; \boldsymbol{\theta}, P) &= \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln p(\mathbf{x}_i, j; \boldsymbol{\theta}_j) \\ &= \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \boldsymbol{\theta}_j) P_j) \end{aligned}$$

When  $\ln p(X; \boldsymbol{\theta}, P)$  is **maximized**?

When **large**  $P(j|\mathbf{x}_i)$ 's are **multiplied** by **large**  $\ln p(\mathbf{x}_i, j; \boldsymbol{\theta}_j)$  's.

# Probabilistic CFO clustering algorithms

For **fixed  $\theta_j$ 's**: Use the Bayes rule  $P(j|\mathbf{x}) = \frac{p(\mathbf{x}|j;\theta_j)P_j}{p(\mathbf{x};\boldsymbol{\theta})}$

For **fixed  $P(j|\mathbf{x})$ 's**: Solve the following maximization problem

$$\begin{aligned} \max_{\theta, P} \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \theta_j) P_j) \\ = \max_{\theta, P} \left[ \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \theta_j)) + \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln P_j \right] \end{aligned}$$

**Subject to the constraint  $\sum_{j=1}^m P_j = 1$ .**

# Mixture models – Expectation-Maximization (EM) algorithm

For **fixed  $\theta_j$ 's**: Use the Bayes rule  $P(j|\mathbf{x}) = \frac{p(\mathbf{x}|j;\theta_j)P_j}{p(\mathbf{x};\Theta)}$

For **fixed  $P(j|\mathbf{x})$ 's**: Solve the following maximization problem

$$\begin{aligned} & \max_{\Theta, P} \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \theta_j)P_j) = \\ \max_{\Theta} & \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \theta_j)) + \max_P \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln P_j \\ & = \max_{\Theta} \sum_{j=1}^m \sum_{i=1}^N P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \theta_j)) + \max_P \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln P_j \end{aligned}$$

**Subject to the constraint**  $\sum_{j=1}^m P_j = 1$ .

The above maximization problem is equivalent to the following maximization **sub-problems**

$$- \theta_j = \operatorname{argmax}_{\theta_j} \sum_{i=1}^N P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \theta_j)), j = 1, \dots, m$$

$$- P \equiv \{P_1, P_2, \dots, P_m\} = \operatorname{argmax}_P \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln P_j, \text{ s.t. } \sum_{j=1}^m P_j = 1 \Leftrightarrow$$

$$P_j = \frac{1}{N} \sum_{i=1}^N P(j|\mathbf{x}_i), j = 1, \dots, m$$

# Probabilistic CFO clustering algorithms

## Generalized probabilistic Algorithmic Scheme (GPrAS)

- Choose  $\theta_j(0), P_j(0)$  as **initial estimates** for  $\theta_j, P_j$ , respectively,  $j = 1, \dots, m$

- $t=0$

- **Repeat**

– For  $i=1$  to  $N$  % *Expectation step*

o For  $j=1$  to  $m$

$$P(j|\mathbf{x}_i; \Theta^{(t)}, P^{(t)}) = \frac{p(x_i|j; \theta_j^{(t)})P_j^{(t)}}{\sum_{q=1}^m p(x_i|q; \theta_q^{(t)})P_q^{(t)}} \equiv \gamma_{ji}^{(t)}$$

o End {For- $j$ }

– End {For- $i$ }

–  $t=t+1$

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

o Set

$$\theta_j^{(t)} = \operatorname{argmax}_{\theta_j} \sum_{i=1}^N \gamma_{ji}^{(t-1)} \ln \left( p(\mathbf{x}_i|j; \theta_j) \right), j = 1, \dots, m$$

$$P_j^{(t)} = \frac{1}{N} \sum_{i=1}^N \gamma_{ji}^{(t-1)}, j = 1, \dots, m$$

– End {For- $j$ }

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



# Probabilistic CFO clustering algorithms

**Remark:** The above algorithm is an instance of the more general **Expectation-Maximization (EM)** framework.

*GPrAS – The case of normal pdfs*

Each **cluster** is **modeled** by a **normal distribution**

$$p(\mathbf{x}|j; \mu_j, \Sigma_j) = \frac{1}{(2\pi)^l |\Sigma_j|^{1/2}} \exp\left(-\frac{(\mathbf{x} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x} - \mu_j)}{2}\right), j = 1, \dots, m$$

In this case  $\theta_j = \{\mu_j, \Sigma_j\}$ .

$$\{\mu_j, \Sigma_j\} = \operatorname{argmax}_{\{\mu_j, \Sigma_j\}} \sum_{i=1}^N P(j|\mathbf{x}_i) \ln\left(p(\mathbf{x}_i|j; \mu_j, \Sigma_j)\right)$$

Equating the gradient of the above function wrt  $\mu_j, \Sigma_j$  to  $\mathbf{0}$  and  $O$ , respectively, we have

$$\mu_j = \frac{\sum_{i=1}^N P(j|\mathbf{x}_i) \mathbf{x}_i}{\sum_{i=1}^N P(j|\mathbf{x}_i)}$$

$$\Sigma_j = \frac{\sum_{i=1}^N P(j|\mathbf{x}_i) (\mathbf{x}_i - \mu_j) (\mathbf{x}_i - \mu_j)^T}{\sum_{i=1}^N P(j|\mathbf{x}_i)}$$

# Probabilistic CFO clustering algorithms

## GPrAS – The normal pdfs case

- Choose  $\mu_j(0), \Sigma_j(0), P_j(0)$  as **initial estimates** for  $\mu_j, \Sigma_j, P_j$ , resp.,  $j = 1, \dots, m$

- $t=0$

- **Repeat**

- For  $i=1$  to  $N$  % *Expectation step*

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

$$P(j|\mathbf{x}_i; \Theta^{(t)}, P^{(t)}) = \frac{p(\mathbf{x}_i|j; \theta_j^{(t)})P_j^{(t)}}{\sum_{q=1}^m p(\mathbf{x}_i|q; \theta_q^{(t)})P_q^{(t)}} \equiv \gamma_{ji}^{(t)}$$

- o End {For- $j$ }

- End {For- $i$ }

- $t=t+1$

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

- o Set

$$\mu_j^{(t)} = \frac{\sum_{i=1}^N \gamma_{ji}^{(t-1)} \mathbf{x}_i}{\sum_{i=1}^N \gamma_{ji}^{(t-1)}}, \quad \Sigma_j^{(t)} = \frac{\sum_{i=1}^N \gamma_{ji}^{(t-1)} (\mathbf{x}_i - \mu_j) (\mathbf{x}_i - \mu_j)^T}{\sum_{i=1}^N \gamma_{ji}^{(t-1)}} \quad j = 1, \dots, m$$

$$P_j^{(t)} = \frac{1}{N} \sum_{i=1}^N \gamma_{ji}^{(t-1)}, \quad j = 1, \dots, m$$

- End {For- $j$ }

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

# Probabilistic CFO clustering algorithms

## GPrAS – The normal pdfs case

- Choose  $\boldsymbol{\mu}_j(0), \boldsymbol{\Sigma}_j(0), P_j(0)$  as **initial estimates** for  $\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j, P_j$ , resp.,  $j = 1, \dots, m$

- $t=0$

- **Repeat**

- For  $i=1$  to  $N$  *% Expectation step*

- $P(C_j|\mathbf{x}; \Theta(t))$

$$= \frac{|\boldsymbol{\Sigma}_j(t)|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j(t))^T \boldsymbol{\Sigma}_j^{-1}(t)(\mathbf{x} - \boldsymbol{\mu}_j(t))\right) P_j(t)}{\sum_{k=1}^m |\boldsymbol{\Sigma}_k(t)|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k(t))^T \boldsymbol{\Sigma}_k^{-1}(t)(\mathbf{x} - \boldsymbol{\mu}_k(t))\right) P_k(t)}$$

- o End {For-j}

- End {For-i}

- $t=t+1$

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

- o Set

$$\boldsymbol{\mu}_j^{(t)} = \frac{\sum_{i=1}^N \gamma_{ji}^{(t-1)} \mathbf{x}_i}{\sum_{i=1}^N \gamma_{ji}^{(t-1)}}, \quad \boldsymbol{\Sigma}_j^{(t)} = \frac{\sum_{i=1}^N \gamma_{ji}^{(t-1)} (\mathbf{x}_i - \boldsymbol{\mu}_j) (\mathbf{x}_i - \boldsymbol{\mu}_j)^T}{\sum_{i=1}^N \gamma_{ji}^{(t-1)}} \quad j = 1, \dots, m$$

$$P_j^{(t)} = \frac{1}{N} \sum_{i=1}^N \gamma_{ji}^{(t-1)}, \quad j = 1, \dots, m$$

- End {For-j}

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

# Probabilistic CFO clustering algorithms

## Remark:

- The above scheme is **more computationally demanding** since it requires the inversion of the  $m$  covariance matrices at each iteration step. Two ways to deal with this problem are:
  - The use of a **single covariance matrix for all clusters**.
  - The use of **different diagonal covariance matrices**.

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

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

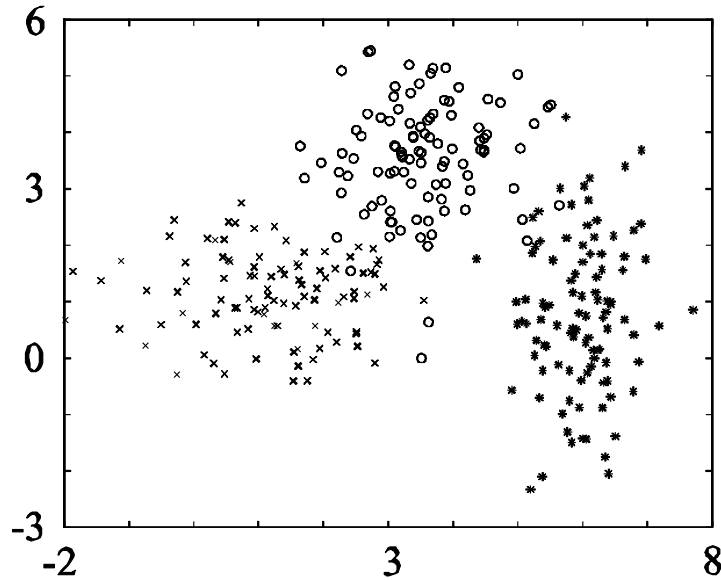
and covariance matrices

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

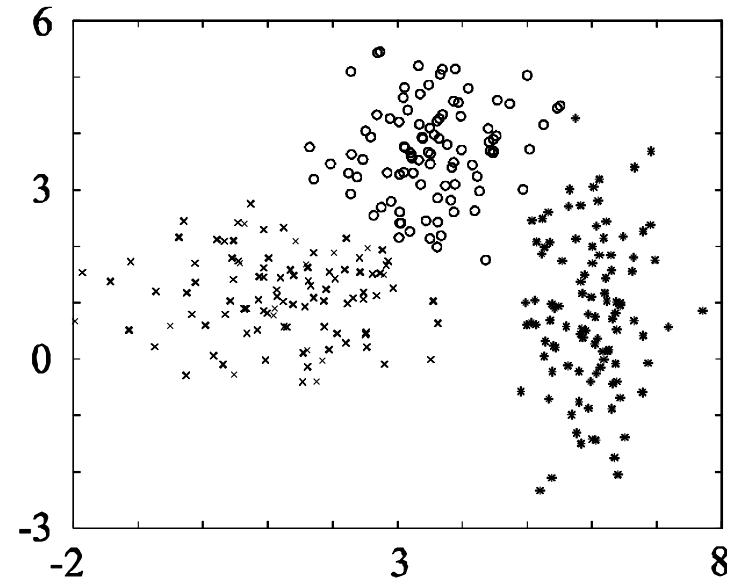
respectively.

A group of 100 vectors stem from each distribution. These form the data set  $X$ .

# Probabilistic CFO clustering algorithms



(a) The data set



(b) Results of GMDAS

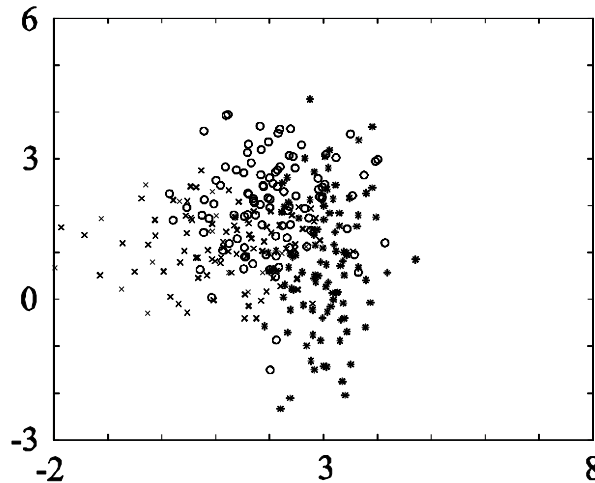
Confusion matrix:

	<i>Cluster 1</i>	<i>Cluster 2</i>	<i>Cluster 3</i>
<i>1<sup>st</sup> distribution</i>	99	0	1
<i>2<sup>nd</sup> distribution</i>	0	100	0
<i>3<sup>rd</sup> distribution</i>	3	4	93

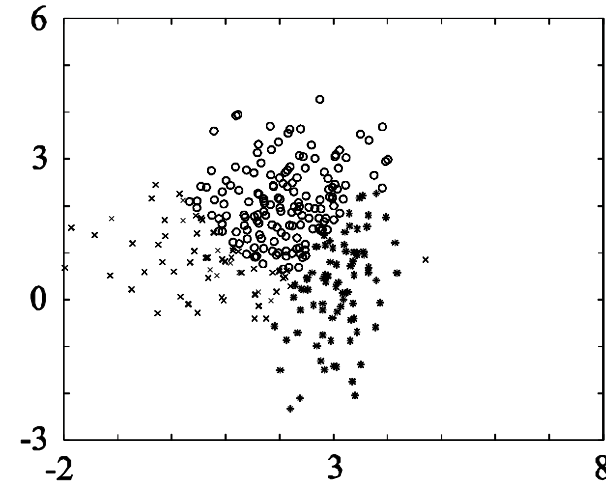
The algorithm reveals accurately the underlying structure.

# Probabilistic CFO clustering algorithms

(b) The same as (a) but now  $\underline{\mu}_1=[1, 1]^T$ ,  $\underline{\mu}_2=[2, 2]^T$ ,  $\underline{\mu}_3=[3, 1]^T$  (The clusters are closer).



(a)  
The data set



(b)  
Results of GMDAS

Confusion matrix:

	<i>Cluster 1</i>	<i>Cluster 2</i>	<i>Cluster 3</i>
<i>1<sup>st</sup> distribution</i>	85	4	11
<i>2<sup>nd</sup> distribution</i>	35	56	9
<i>3<sup>rd</sup> distribution</i>	26	0	74

The algorithm reveals the underlying structure less accurately.

# Probabilistic CFO clustering algorithms

Example  $\mathbf{x}_1 = [0 \ 0]^T$ ,  $\mathbf{x}_2 = [3 \ 0]^T$ ,  $\mathbf{x}_3 = [0 \ 3]^T$ ,  $\mathbf{x}_4 = [12 \ 12]^T$ ,  $\mathbf{x}_5 = [15 \ 12]^T$ ,  $\mathbf{x}_6 = [12 \ 15]^T$

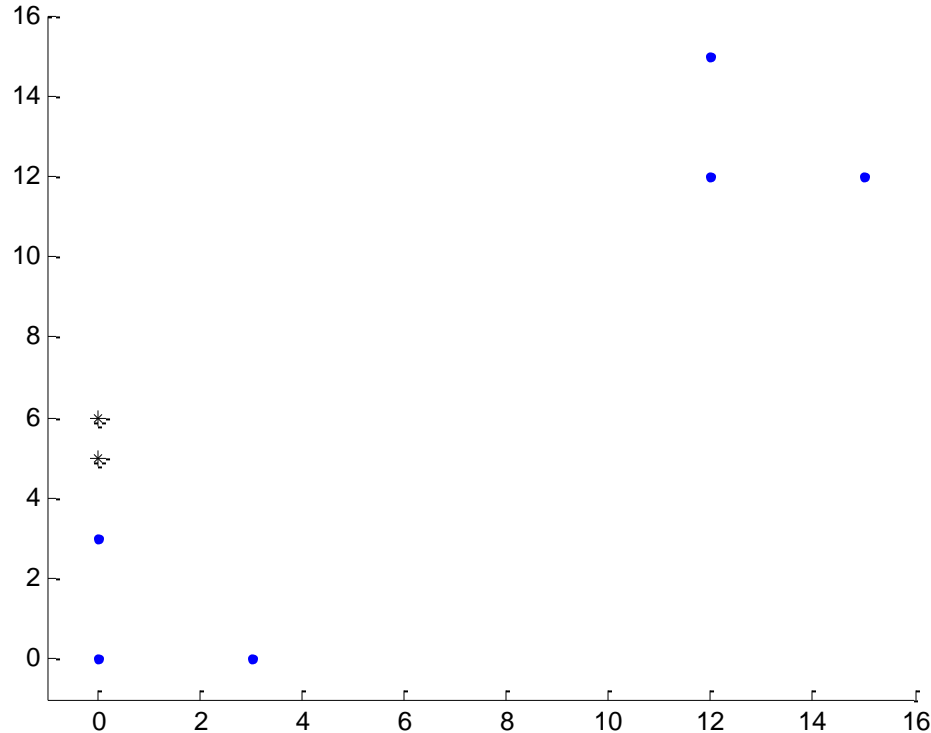
**Initially:**

$$\theta_1(0) = [0, 5]^T$$

$$\theta_2(0) = [0, 6]^T$$

$$P_1(0) = 0.1$$

$$P_2(0) = 0.9$$



$$p(\mathbf{x}|1) = \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_1\|^2), \quad P(1|\mathbf{x}) = \frac{p(\mathbf{x}|1)P_1}{p(\mathbf{x})}$$

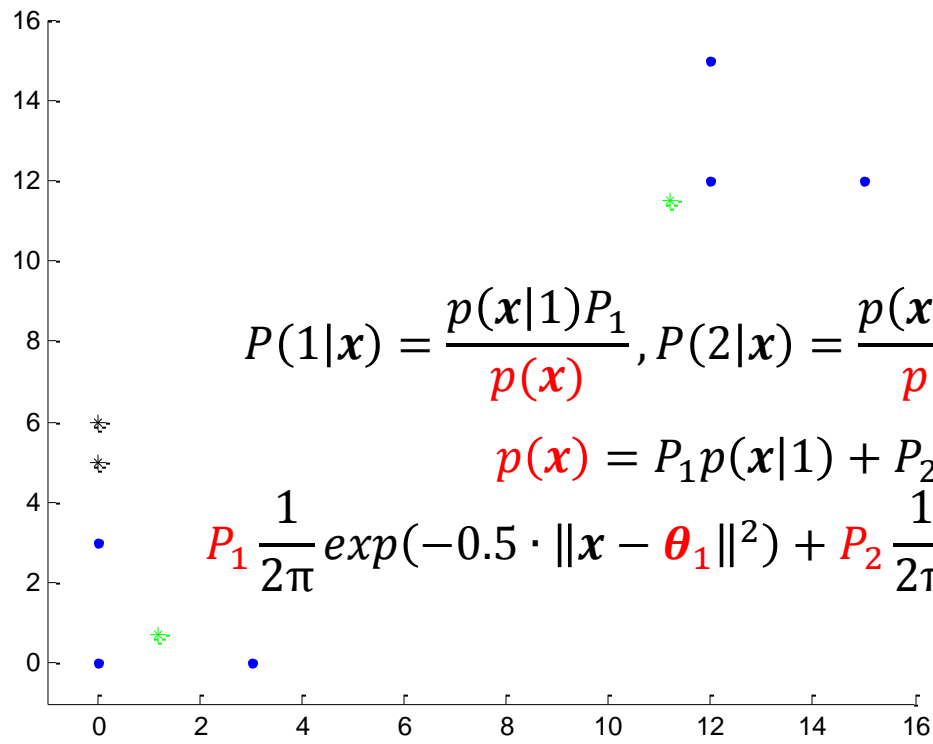
$$p(\mathbf{x}|2) = \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_2\|^2), \quad P(2|\mathbf{x}) = \frac{p(\mathbf{x}|2)P_2}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = P_1 p(\mathbf{x}|1) + P_2 p(\mathbf{x}|2) = P_1 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_1\|^2) + P_2 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_2\|^2)$$

$$\ln p(X; \theta, P) = \sum_{i=1}^N [P(1|\mathbf{x}_i) \ln(p(\mathbf{x}_i|1; \boldsymbol{\theta}_1)P_1) + P(2|\mathbf{x}_i) \ln(p(\mathbf{x}_i|2; \boldsymbol{\theta}_2)P_2)]$$

# Probabilistic CFO clustering algorithms

Example  $\mathbf{x}_1 = [0 \ 0]^T$ ,  $\mathbf{x}_2 = [3 \ 0]^T$ ,  $\mathbf{x}_3 = [0 \ 3]^T$ ,  $\mathbf{x}_4 = [12 \ 12]^T$ ,  $\mathbf{x}_5 = [15 \ 12]^T$ ,  $\mathbf{x}_6 = [12 \ 15]^T$



$$P(1|\mathbf{x}) = \frac{p(\mathbf{x}|1)P_1}{p(\mathbf{x})}, P(2|\mathbf{x}) = \frac{p(\mathbf{x}|2)P_2}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = P_1 p(\mathbf{x}|1) + P_2 p(\mathbf{x}|2) =$$

$$P_1 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_1\|^2) + P_2 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_2\|^2)$$

1<sup>st</sup> iteration:

A posteriori probs

	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$	$\mathbf{x}_5$	$\mathbf{x}_6$
$P(1 \mathbf{x})$	0.9645	0.9645	0.5751	0.0002	0.0002	0.0000
$P(2 \mathbf{x})$	0.0355	0.0355	0.4249	0.9998	0.9998	1.0000

$$\boldsymbol{\theta}_1(1) = [1.1572 \ 0.6906]^T$$

$$\boldsymbol{\theta}_2(1) = [11.1864 \ 11.5207]^T$$

$$P_1(1) = 0.4174$$

$$P_2(1) = 0.5826$$



# Probabilistic CFO clustering algorithms

Example  $x_1 = [0 \ 0]^T$ ,  $x_2 = [3 \ 0]^T$ ,  $x_3 = [0 \ 3]^T$ ,  $x_4 = [12 \ 12]^T$ ,  $x_5 = [15 \ 12]^T$ ,  $x_6 = [12 \ 15]^T$

1<sup>st</sup> iteration (in more detail):

A. Expectation step - A posteriori probs

$$\begin{aligned}
 P(1|x_1) &= \frac{P_1(0) \frac{1}{2\pi} \exp(-0.5 \cdot \|x_1 - \theta_1(0)\|^2)}{P_1(0) \frac{1}{2\pi} \exp(-0.5 \cdot \|x_1 - \theta_1(0)\|^2) + P_2(0) \frac{1}{2\pi} \exp(-0.5 \cdot \|x_1 - \theta_2(0)\|^2)} = \\
 &= \frac{0.1 \frac{1}{2\pi} \exp\left(-0.5 \cdot \left\| \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 5 \end{bmatrix} \right\|^2\right)}{0.1 \frac{1}{2\pi} \exp\left(-0.5 \cdot \left\| \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 5 \end{bmatrix} \right\|^2\right) + 0.9 \frac{1}{2\pi} \exp\left(-0.5 \cdot \left\| \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 6 \end{bmatrix} \right\|^2\right)} = \\
 &= \frac{0.1 \frac{1}{2\pi} \exp(-12.5)}{0.1 \frac{1}{2\pi} \exp(-12.5) + 0.9 \frac{1}{2\pi} \exp(-18)} = \mathbf{0.9645}
 \end{aligned}$$

In a similar way we have

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
$P(1 x)$	0.9645	0.9645	0.5751	0.0002	0.0002	0.0000
$P(2 x)$	0.0355	0.0355	0.4249	0.9998	0.9998	1.0000

# Probabilistic CFO clustering algorithms

Example  $x_1 = [0 \ 0]^T$ ,  $x_2 = [3 \ 0]^T$ ,  $x_3 = [0 \ 3]^T$ ,  $x_4 = [12 \ 12]^T$ ,  $x_5 = [15 \ 12]^T$ ,  $x_6 = [12 \ 15]^T$

1<sup>st</sup> iteration (in more detail):

B. Maximization step - Parameters  $\theta_1, \theta_2, P_1, P_2$

$$\theta_1(1) = \frac{0.9645 \cdot x_1 + 0.9645 \cdot x_2 + 0.5751 \cdot x_3 + 0.0002 \cdot x_4 + 0.0002 \cdot x_5 + 0.0000 \cdot x_6}{0.9645 + 0.9645 + 0.5751 + 0.0002 + 0.0002 + 0.0000} =$$

$$\frac{0.9645 \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} + 0.9645 \cdot \begin{bmatrix} 3 \\ 0 \end{bmatrix} + 0.5751 \cdot \begin{bmatrix} 0 \\ 3 \end{bmatrix} + 0.0002 \cdot \begin{bmatrix} 12 \\ 12 \end{bmatrix} + 0.0002 \cdot \begin{bmatrix} 15 \\ 12 \end{bmatrix} + 0.0000 \cdot \begin{bmatrix} 12 \\ 15 \end{bmatrix}}{0.9645 + 0.9645 + 0.5751 + 0.0002 + 0.0002 + 0.0000} =$$

$$= \begin{bmatrix} 1.1572 \\ 0.6906 \end{bmatrix}$$

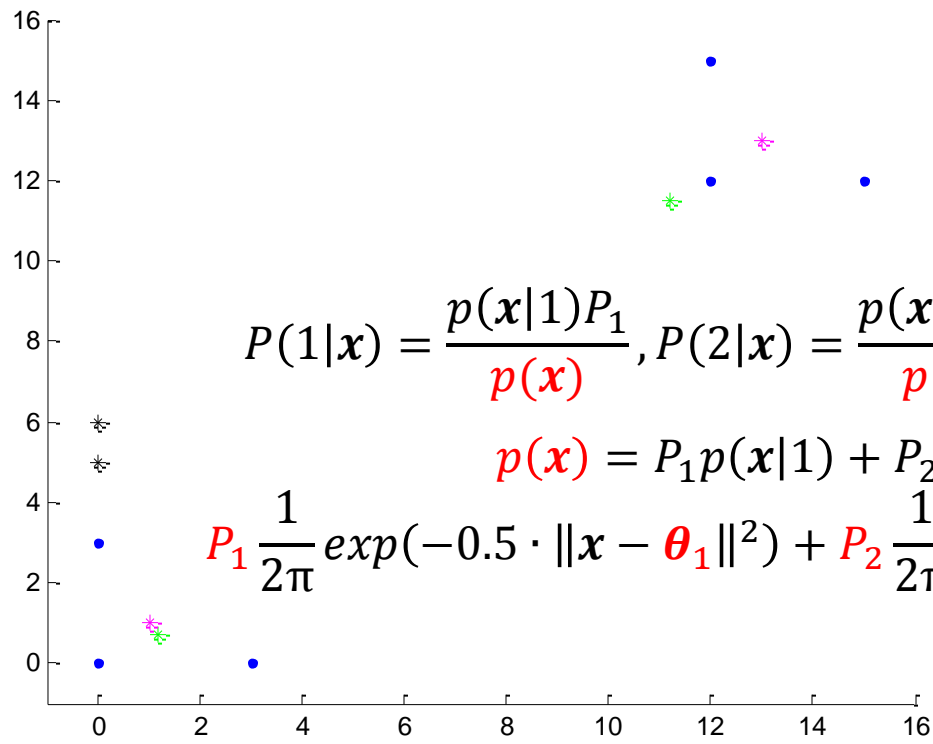
In a similar way we have  $\theta_2(1) = \begin{bmatrix} 11.1864 \\ 11.5207 \end{bmatrix}$

$$P_1(1) = \frac{0.9645 + 0.9645 + 0.5751 + 0.0002 + 0.0002 + 0.0000}{6} = 0.4174$$

In a similar way we have  $P_2(1) = 0.5826$

# Probabilistic CFO clustering algorithms

Example  $x_1 = [0 \ 0]^T$ ,  $x_2 = [3 \ 0]^T$ ,  $x_3 = [0 \ 3]^T$ ,  $x_4 = [12 \ 12]^T$ ,  $x_5 = [15 \ 12]^T$ ,  $x_6 = [12 \ 15]^T$



$$P(1|\mathbf{x}) = \frac{p(\mathbf{x}|1)P_1}{p(\mathbf{x})}, P(2|\mathbf{x}) = \frac{p(\mathbf{x}|2)P_2}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = P_1 p(\mathbf{x}|1) + P_2 p(\mathbf{x}|2) =$$

$$P_1 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_1\|^2) + P_2 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_2\|^2)$$

**2<sup>nd</sup> iteration:**  
**A posteriori probs**

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
$P(1 \mathbf{x})$	1.0000	1.0000	1.0000	0.0000	0.0000	0.0000
$P(2 \mathbf{x})$	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000

$$\boldsymbol{\theta}_1(2) = [1 \ 1]^T$$

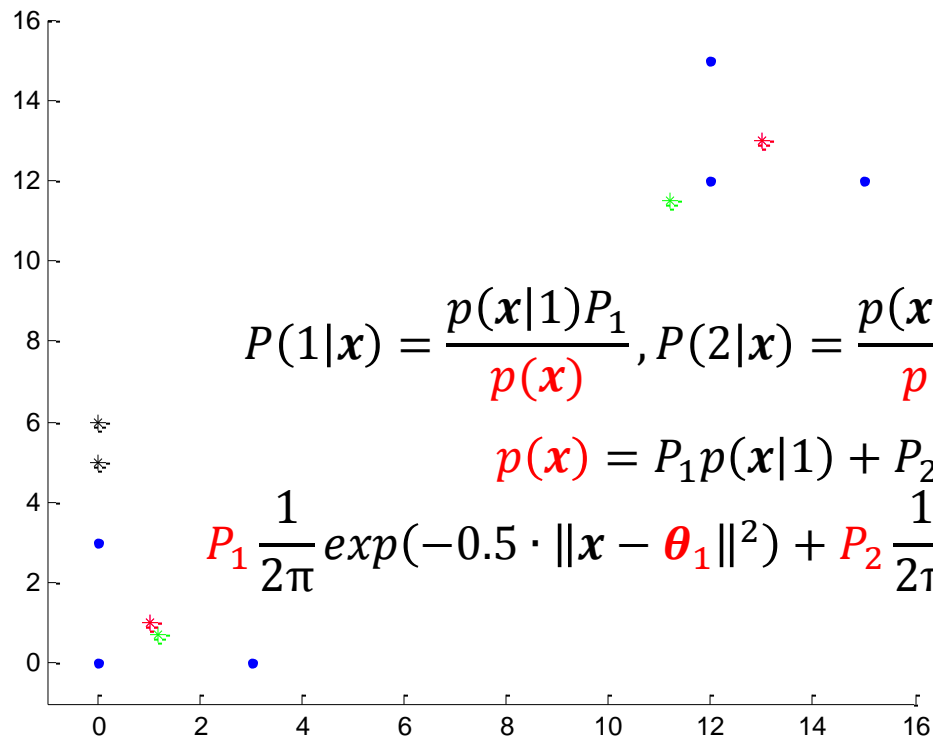
$$\boldsymbol{\theta}_2(2) = [13 \ 13]^T$$

$$P_1(2) = 0.5$$

$$P_2(2) = 0.5$$

# Probabilistic CFO clustering algorithms

Example  $\mathbf{x}_1 = [0\ 0]^T$ ,  $\mathbf{x}_2 = [3\ 0]^T$ ,  $\mathbf{x}_3 = [0\ 3]^T$ ,  $\mathbf{x}_4 = [12\ 12]^T$ ,  $\mathbf{x}_5 = [15\ 12]^T$ ,  $\mathbf{x}_6 = [12\ 15]^T$



$$P(1|\mathbf{x}) = \frac{p(\mathbf{x}|1)P_1}{p(\mathbf{x})}, P(2|\mathbf{x}) = \frac{p(\mathbf{x}|2)P_2}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = P_1 p(\mathbf{x}|1) + P_2 p(\mathbf{x}|2) =$$

$$P_1 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_1\|^2) + P_2 \frac{1}{2\pi} \exp(-0.5 \cdot \|\mathbf{x} - \boldsymbol{\theta}_2\|^2)$$

**3<sup>rd</sup> iteration:**  
**A posteriori probs**

	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$	$\mathbf{x}_5$	$\mathbf{x}_6$
$P(1 \mathbf{x})$	1.0000	1.0000	1.0000	0.0000	0.0000	0.0000
$P(2 \mathbf{x})$	0.0000	0.0000	0.0000	1.0000	1.0000	1.0000

$$\boldsymbol{\theta}_1(3) = [1\ 1]^T$$

$$\boldsymbol{\theta}_2(3) = [13\ 13]^T$$

$$P_1(3) = 0.5$$

$$P_2(3) = 0.5$$