# Finite Mixture Models and Clustering

## Mohamed Nadif

LIPADE, Université Paris Descartes, France

## Outline

| Types | Algorithms | Criteria | dissimilarity/similarity measures |
|---|---|---|---|
| Continuous | k-means | $\sum_{i,k} z_{ik} D(\mathbf{x}_i, \boldsymbol{\mu}_k)$ $\mathbf{x}_i, \boldsymbol{\mu}_k \in \mathbb{R}^p$ | $D(\mathbf{x}_i, \boldsymbol{\mu}_k) = \sum_j (x_{ij} - \mu_{kj})^2$ |
| Contingency | k-means-$\chi^2$ | $\sum_{i,k} z_{ik} D_{\chi^2}(\mathbf{x}_i, \boldsymbol{\mu}_k)$ $\mathbf{x}_i = (\frac{x_{i1}}{x_{i.}}, \ldots, \frac{x_{ip}}{x_{i.}})^T$ $\mathbf{x}_i, \boldsymbol{\mu}_k \in [0,1]^p$ | $D_{\chi^2}(\mathbf{x}_i, \boldsymbol{\mu}_k) = \sum_j \frac{1}{x_{.j}} (\frac{x_{ij}}{x_{i.}} - \mu_{kj})^2$ |
| Binary | k-modes | $\sum_{i,k} z_{ik} D(\mathbf{x}_i, \mathbf{a}_k)$ $\mathbf{x}_i, \mathbf{a}_k \in \{0,1\}^p$ | $D(\mathbf{x}_i, \mathbf{a}_k) = \sum_j |x_{ij} - a_{kj}|$ |
| Categorical | k-modes | $\sum_{i,k} z_{ik} D(\mathbf{x}_i, \boldsymbol{\lambda}_k)$ $\mathbf{x}_i, \boldsymbol{\lambda}_k \in \{1, \ldots, m^j\}^p$ | $D(\mathbf{x}_i, \mathbf{a}_k) = \sum_j \delta(x_{ij}, \lambda_{kj})$ $\delta(x_{ij}, \lambda_{kj}) = 1$ if $x_{ij} = \lambda_{kj}$ $\delta(x_{ij}, \lambda_{kj}) = 0$ if $x_{ij} \neq \lambda_{kj}$ |
| Directional | Sk-means | $\sum_{i,k} z_{ik} \cos(\mathbf{x}_i, \boldsymbol{\mu}_k)$ $\mathbf{x}_i, \boldsymbol{\mu}_k \in [0,1]^p$ | $\cos(\mathbf{x}_i, \boldsymbol{\mu}_k)$ |

**Table:** Criteria and algorithms

### Classical clustering methods

- Clustering methods hierarchical and nonhierarchical methods have advantages and disadvantages
- Disadvantages. They are for the most part heuristic techniques derived from empirical methods
- Difficulties to take into account the characteristics of clusters (shapes, proportions, volume etc.)
- Geometrical approach: Clustering with "adaptives" distances: $d_{M_k}(\mathbf{x}, \mathbf{y}) = ||\mathbf{x} - \mathbf{y}||_{M_k}$
- In fact, the principal question "does it exist a model ?"

### Mixture Approach

- MA have attracted much attention since 1990.
- It is undoubtedly a very useful contribution to clustering
    1. It offers considerable flexibility
    2. provides solutions to the problem of the number of clusters
    3. Its associated estimators of posterior probabilities give rise to a fuzzy or hard clustering using the a MAP
    4. It permits to give a meaning to certain classical criteria
- Finite Mixture Models by (McLachlan and Peel, 2000)
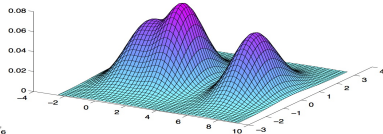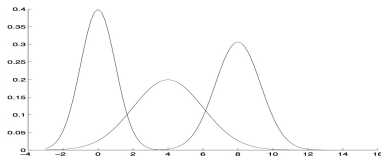
## Outline

### Definition of the model

- In model-based clustering it is assumed that the data are generated by a mixture of underlying probability distributions, where each component $k$ of the mixture represents a cluster. Thus, the data matrix is assumed to be an i.i.d sample $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$ where $\boldsymbol{x}_i = (x_{i1}, \ldots, x_{ip}) \in \mathbb{R}^p$ from a probability distribution with density

$$f(\boldsymbol{x}_i; \Theta) = \sum_{k=1}^{g} \pi_k \varphi(\boldsymbol{x}_i; \boldsymbol{\alpha}_k),$$



where

- $\varphi(\,.\;; \boldsymbol{\alpha}_k)$ is the density of an observation $\boldsymbol{x}_i$ from the $k$-th component
- $\boldsymbol{\alpha}_k$'s are the corresponding class parameters. These densities belong to the same parametric family
- The parameter $\pi_k$ is the probability that an object belongs to the $k$-th component
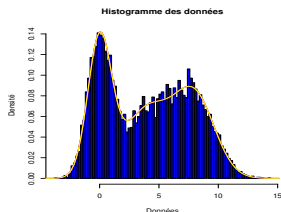
## Gaussian mixture model in $\mathbb{R}^1$

- n=9000, p=1, g=3
- $\varphi(., \alpha_k)$ a Gaussian density $\alpha_k = (m_k, s_k)$
- $\pi_1 = \pi_2 = \pi_3 = \frac{1}{3}$

The mixture density of the observed data x can be written as

$$f(\mathbf{X}; \Theta) = \prod_{i=1}^{n} \sum_{k=1}^{g} \pi_k \frac{1}{s_k \sqrt{2\pi}} \exp(-\frac{1}{2}(\frac{x_i - m_k}{s_k})^2)$$

## Mixture of 3 densities

## Gaussian mixture model in $\mathbb{R}^2$: N((2,1);1/3) and N((7,5);1)

```
X=matrix(nrow=1000,ncol=2)
for (i in 1:1000)
{
Z=rbinom(1,1,2/3)
if (Z==1){
X[i,1]=rnorm(1,2,1)
X[i,2]=rnorm(1,1,1)
}
else
{
X[i,1]=rnorm(1,7,1)
X[i,2]=rnorm(1,5,1)
}
}
plot(X)
```

**Likelihood of observed data X**

- The parameter of this model is the vector $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\alpha})$ containing the mixing proportions $\boldsymbol{\pi} = (\pi_1, ..., \pi_g)$ and the vector $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1, ..., \boldsymbol{\alpha}_g)$ of parameters of each component. The mixture density of the observed data **X** can be expressed as

$$f(\mathbf{X}; \Theta) = \prod_{i=1}^{n} \sum_{k=1}^{g} \pi_k \varphi(\boldsymbol{x}_i; \boldsymbol{\alpha}_k).$$

**Bernoulli mixture model**

- For instance, for binary data with $\boldsymbol{x}_i \in \{0, 1\}^p$, using multivariate Bernoulli distributions for each component, the mixture density of the observed data **x** can be written as

$$f(\mathbf{X}; \Theta) = \prod_{i=1}^{n} \sum_{k=1}^{g} \pi_k \prod_{j=1}^{p} \alpha_{kj}^{x_{ij}} (1 - \alpha_{kj})^{1 - x_{ij}}$$

where $\Theta = \{\pi_1, \ldots, \pi_g, \boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}_g\}$ with $\boldsymbol{\alpha}_k = (\alpha_{k1}, \ldots, \alpha_{kp})$ and $\alpha_{kj} \in [0, 1]$

### ML and CML approaches

- The problem of clustering can be studied in the mixture model using two different approaches: the maximum likelihood approach (ML) and the classification likelihood approach (CML)

  **1** The ML approach (Day, 1969): It estimates the parameters of the mixture, and the partition on the objects is derived from these parameters using the maximum a posteriori principle (MAP). The maximum likelihood estimation of the parameters results in an optimization of the log-likelihood of the observed sample

  $$L_M(\Theta) = L(\Theta; \mathbf{X}) = \sum_{i=1}^{n} \log \left( \sum_{k=1}^{g} \pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k) \right)$$

  **2** The CML approach (Symons, 1981): It estimates the parameters of the mixture and the partition *simultaneously* by optimizing the classification log-likelihood

  $$L_C(\mathbf{z}; \Theta) = L(\Theta; \mathbf{X}, \mathbf{z}) = \log f(\mathbf{X}, \mathbf{z}; \Theta) = \sum_{i=1}^{n} \sum_{k=1}^{g} z_{ik} \log \left( \pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k) \right)$$

  or

  $$L_C(\mathbf{z}; \Theta) = \sum_{i=1}^{n} \sum_{k=1}^{g} z_{ik} \log \left( \pi_k \right) + \sum_{i=1}^{n} \sum_{k=1}^{g} z_{ik} \log \left( \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k) \right)$$

## Outline

**Introduction of EM**

- Much effort has been devoted to the estimation of parameters for the mixture model
- Pearson used the method of moments to estimate $\Theta = (m_1, m_2, s_1^2, s_2^2, \pi)$ of a unidimensional Gaussian mixture model with two components

$$f(\mathbf{x}_i; \boldsymbol{\theta}) = \pi\varphi(\mathbf{x}_i; m_1, s_1^2) + (1 - \pi)\varphi(\mathbf{x}_i; m_2, s_2^2)$$

  required to solve polynomial equations of degree nine

- Generally, the appropriate method used in this context is the EM algorithm (Dempster et al., 1977). Two steps Expectation and Maximization
- This algorithm can be applied in different contexts where the model depends on unobserved latent variables. In mixture context $\mathbf{z}$ represents this variable. It denotes which $\mathbf{x}_i$ is from. Then we note $\mathbf{Y} = (\mathbf{X}, \mathbf{z})$ the complete data.
- Starting from the relation between the densities

$$f(\mathbf{Y}, \Theta) = f((\mathbf{X}, \mathbf{z}); \Theta) = f(\mathbf{Y}|\mathbf{X}; \Theta)f(\mathbf{X}; \Theta)$$

  we have

$$\log(f(\mathbf{X}; \Theta)) = \log(f(\mathbf{Y}, \Theta)) - \log(f(\mathbf{Y}|\mathbf{X}; \Theta))$$

  or

$$L_M(\Theta) = L_C(\mathbf{z}; \Theta) - \log f(\mathbf{Y}|\mathbf{X}; \Theta)$$

**Principle of EM**

- Objective: Maximization of $L_M(\Theta)$
- EM relies on iterative procedure based on the conditional expectation of $L_M(\Theta)$ for a value of the current parameter $\Theta'$

$$L_M(\Theta) = Q(\Theta|\Theta') - H(\Theta|\Theta')$$

where $Q(\Theta|\Theta') = \mathbb{E}(L_C(\mathbf{z}; \Theta|\mathbf{X}, \Theta'))$ and $H(\Theta|\Theta') = \mathbb{E}(\log f(\mathbf{Y}|\mathbf{X}; \Theta)|\mathbf{X}, \Theta')$

- Using the Jensen inequality (Dempster et al;, 1977) for fixed $\Theta'$ we have $\forall \Theta, H(\Theta|\Theta') \leq H(\Theta'|\Theta')$. This inequality can be proved

$$H(\Theta|\Theta') - H(\Theta'|\Theta') = \sum_{\mathbf{z} \in \mathcal{Z}} f(\mathbf{z}|\mathbf{X}; \Theta') \log \frac{f(\mathbf{z}|\mathbf{X}; \Theta)}{f(\mathbf{z}|\mathbf{X}; \Theta')}$$

As $\log(x) \leq x - 1$, we have $\log \frac{f(\mathbf{z}|\mathbf{X};\Theta)}{f(\mathbf{z}|\mathbf{X};\Theta')} \leq \frac{f(\mathbf{z}|\mathbf{X};\Theta)}{f(\mathbf{z}|\mathbf{X};\Theta')} - 1$ then

$$H(\Theta|\Theta') - H(\Theta'|\Theta') \leq \sum_{\mathbf{z} \in \mathcal{Z}} f(\mathbf{z}|\mathbf{X}; \Theta) - \sum_{\mathbf{z} \in \mathcal{Z}} f(\mathbf{z}|\mathbf{X}; \Theta') = 1 - 1 = 0$$

## $Q(\Theta|\Theta')$

- The value $\Theta$ maximizing $Q(\Theta|\Theta')$ satisfies the relation $Q(\Theta|\Theta') \geq Q(\Theta'|\Theta')$ and,

$$L_M(\Theta) = Q(\Theta|\Theta') - H(\Theta|\Theta') \geq Q(\Theta'|\Theta') - H(\Theta'|\Theta') = L_M(\Theta')$$

- In the mixture context

$$Q(\Theta|\Theta') = \mathbb{E}(L_C(\mathbf{z}; \Theta|\mathbf{X}, \Theta')) = \sum_{i,k} \mathbb{E}(z_{ik}|\mathbf{X}, \Theta') \log(\pi_k f(\mathbf{x}_i; \alpha_k))$$

Note that $\mathbb{E}(z_{ik}|\mathbf{X}, \Theta') = p(z_{ik} = 1|\mathbf{X}, \Theta')$
As the conditional distribution of the missing data $\mathbf{z}$ given the observed values :

$$f(\mathbf{z}|\mathbf{X}; \Theta) = \frac{f(\mathbf{X}, \mathbf{z}; \boldsymbol{\theta})}{f(\mathbf{X}; \Theta)} = \frac{f(\mathbf{X}|\mathbf{z}; \Theta) f(\mathbf{z}; \Theta)}{f(\mathbf{X}; \Theta)}$$

we have

$$p(z_{ik} = 1|\mathbf{X}, \Theta') = s_{ik} = \frac{\pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k)}{f(\mathbf{x}_i; \boldsymbol{\theta})} = \frac{\pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k)}{\sum_\ell \pi_\ell \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_\ell)} \propto \pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k)$$

## The steps of EM

- The EM algorithm involves constructing, from an initial $\boldsymbol{\theta}^{(0)}$, the sequence $\boldsymbol{\theta}^{(c)}$ satisfying

$$\Theta^{(c+1)} = \text{argmax } Q(\Theta|\Theta^{(c)})$$

and this sequence causes the criterion $L_M(\Theta)$ to grow. The EM algorithm takes the following form

- Initialize by selecting an initial solution $\Theta^{(0)}$
- Repeat the two steps until convergence

  1. E-step: compute $Q(\Theta|\Theta^{(c)})$. Note that in the mixture case this step reduces to the computation of the conditional probabilities $s_{ik}^{(c)}$

  2. M-step: compute $\Theta^{(c+1)}$ maximizing $Q(\Theta, \Theta^{(c)})$. This leads to $\pi_k^{(c+1)} = \frac{1}{n} \sum_i s_{ik}^{(c+1)}$ and the exact formula for the $\alpha_k^{(c+1)}$ will depend on the involved parametric family of distribution probabilities

## Properties of EM

- Under certain conditions, it has been established that EM always converges to a local likelihood maximum
- Simple to implement and it has good behavior in clustering and estimation contexts
- Slow in some situations

## Example

```
library(mixtools)
attach(faithful)
dim(faithful)
waiting
hist(waiting)
d=density(waiting)
plot(d)
wait1 <- normalmixEM(waiting, lambda = .5, mu = c(50, 60), sigma = 5)

plot(wait1, density = TRUE, cex.axis = 1.4, cex.lab = 1.4, cex.main = 1.8,main2 = "Time between Old Faithful

eruptions", xlab2 = "Minutes")
```

## Mixture of 2 densities



Histogram of waiting



density.default(x = waiting)



Time between Old Faithful eruptions

## An other interpretation of EM

### Hathaway interpretation of EM: classical mixture model context

- EM = alternated maximization of the fuzzy clustering criterion

$$F_C(\mathbf{s}, \Theta) = L_C(\mathbf{s}; \Theta) + H(\mathbf{s})$$

  - $\mathbf{s} = (s_{ik})$: fuzzy partition
  - $L_C(\mathbf{s}, \Theta) = \sum_{i,k} s_{ik} \log(\pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k))$: fuzzy classification log-likelihood
  - $H(\mathbf{s}) = -\sum_{i,k} s_{ik} \log s_{ik}$ : entropy function

### Algorithm

- Maximizing $F_C(\mathbf{s}, \Theta)$ w.r. to $\mathbf{s}$ yields the E-step
- Maximizing $F_C(\mathbf{s}, \Theta)$ w.r. to $\Theta$ yields the M-step

### CEM algorithm

- In the CML approach the partition is added to the parameters to be estimated. The maximum likelihood estimation of these new parameters results in an optimization of the complete data log-likelihood. This optimization can be performed using the following Classification EM (CEM) algorithm (Celeux and Govaert, 1992), a variant of EM, which converts the $s_{ik}$'s to a discrete classification in a C-step before performing the M-step:

  - E-step: compute the posterior probabilities $s_{ik}^{(c)}$.
  - C-step: the partition $\mathbf{z}^{(c+1)}$ is defined by assigning each observation $\mathbf{x}_i$ to the cluster which provides the maximum current posterior probability.
  - M-step: compute the maximum likelihood estimate $(\pi_k^{(c+1)}, \boldsymbol{\alpha}_k^{(c+1)})$ using the $k$-th cluster. This leads to $\pi_k^{(c+1)} = \frac{1}{n} \sum_i z_{ik}^{(c+1)}$ and the exact formula for the $\boldsymbol{\alpha}_k^{(c+1)}$ will depend on the involved parametric family of distribution probabilities

### Properties of CEM

- Simple to implement and it has good practical behavior in clustering context
- Faster than EM and scalable
- Some difficulties when the clusters are not well separated

### Link between CEM and the dynamical clustering methods

| Dynamical clustering method | The CEM algorithm |
|---|---|
| Assignation-step | E-step |
| $z_k = \{i; d(\mathbf{x}_i, \mathbf{a}_k) \leq d(\mathbf{x}_i, \mathbf{a}'_k); k' \neq k\}$ | Compute $s_{ik} \propto \pi_k \varphi(\mathbf{x}_i, \alpha_k)$ |
| | C-step |
| | $z_k = \{i; s_{ik} \geq s_{ik'}; k' \neq k\}$ |
| | $z_k = \{i; -\log(\pi_k \varphi(\mathbf{x}_i, \alpha_k)) \leq -\log(\pi_k \varphi(\mathbf{x}_i, \alpha'_k)); k' \neq k\}$ |
| Representation-step | M-step |
| Compute the center $\mathbf{a}_k$ of each cluster | Compute the $\pi_k$'s and $\alpha_k$ |

### Density and distance

- When the proportions are supposed equal we can propose a *"distance"* or a dissimilarity measure $D$ by taking $\varphi(\mathbf{x}_i, \alpha_k) = \exp(-D(\mathbf{x}_i, \mathbf{a}_k))$ then

$$D(\mathbf{x}_i, \mathbf{a}_k) = -\log(\varphi(\mathbf{x}_i, \alpha_k))$$

and the criterion to optimize is

$$-\sum_i \sum_k z_{ik} D(\mathbf{x}_i, \mathbf{a}_k)$$

### Classical algorithms

- Classical k-means, k-means with chisquare metric, k-modes (k-means for categorical data)

# Outline

**Binary data**

- For binary data, considering the conditional independence model (independence for each component), the mixture density of the observed data x can be written as

$$f(\mathbf{X}; \Theta) = \prod_i \sum_k \pi_k \prod_j \alpha_{kj}^{x_{ij}} (1 - \alpha_{kj})^{1-x_{ij}}$$

  where $x_{ij} \in \{0, 1\}$, $\boldsymbol{\alpha}_k = (\alpha_{k1}, \ldots, \alpha_{kp})$ and $\alpha_{kj} \in [0, 1]$

- Latent Class Model

- The different steps of EM algorithm

  1. E-step: compute $s_{ik}$
  2. M-step: $\alpha_{kj} = \frac{\sum_i s_{ik} x_{ij}}{\sum_i s_{ik}}$ and $\pi_k = \frac{\sum_i s_{ik}}{n}$

- The different steps of CEM algorithm

  1. E-step: compute $s_{ik}$
  2. C-step: compute **z**
  3. M-step: $\alpha_{kj} = \frac{\sum_i z_{ik} x_{ij}}{\sum_i z_{ik}} = \%1$ and $\pi_k = \frac{\#z_k}{n}$

**Parsimonious model**

Several parsimonious models can be proposed by imposing constraints s on the parameters

$$f(\mathbf{x}_i; \Theta) = \sum_k \pi_k \prod_j \varepsilon_{kj}^{|x_{ij} - a_{kj}|} (1 - \varepsilon_{kj})^{1 - |x_{ij} - a_{kj}|}$$

where

$$\begin{cases} a_{kj} = 0, \varepsilon_{kj} = \alpha_{kj} & \text{if } \alpha_{kj} < 0.5 \\ a_{kj} = 1, \varepsilon_{kj} = 1 - \alpha_{kj} & \text{if } \alpha_{kj} > 0.5 \end{cases}$$

- The parameter $\boldsymbol{\alpha}_k$ is replaced by the two parameters $\mathbf{a}_k$ and $\varepsilon_k$

   Example: $\boldsymbol{\alpha}_k = (0.7, 0.3, 0.4, 0.6)$ then $\mathbf{a}_k = (1, 0, 0, 1)$ and $\varepsilon_k = (0.3, 0.3, 0.4, 0.4)$
   - The binary vector $a_k$ represents the center of the cluster $z_k$, each $a_{kj}$ indicates the most frequent binary value
   - The binary vector $\varepsilon_k \in ]0, 1/2[^p$ represents the degrees of heterogeneity of the cluster $z_k$, each $\varepsilon_{kj}$ represents the probability of $j$ to have the value different from that of the center,
     - $p(x_{ij} = 1 | a_{kj} = 0) = p(x_{ij} = 0 | a_{kj} = 1) = \varepsilon_{kj}$
     - $p(x_{ij} = 0 | a_{kj} = 0) = p(x_{ij} = 1 | a_{kj} = 1) = 1 - \varepsilon_{kj}$

- 8 Models assuming proportions equal or not : $[\varepsilon_{kj}]$, $[\varepsilon_k]$, $\varepsilon_j$, $[\varepsilon]$

**Binary data matrix and reorganized data matrix**

|    | a | b | c | d | e |    | a | b | c | d | e |
|----|---|---|---|---|---|----|---|---|---|---|---|
| 1  | 1 | 0 | 1 | 0 | 1 | 1  | 1 | 0 | 1 | 0 | 1 |
| 2  | 0 | 1 | 0 | 1 | 0 | 4  | 1 | 0 | 1 | 0 | 0 |
| 3  | 1 | 0 | 0 | 0 | 0 | 8  | 1 | 0 | 1 | 0 | 1 |
| 4  | 1 | 0 | 1 | 0 | 0 | 2  | 0 | 1 | 0 | 1 | 0 |
| 5  | 0 | 1 | 0 | 1 | 1 | 5  | 0 | 1 | 0 | 1 | 1 |
| 6  | 0 | 1 | 0 | 0 | 1 | 6  | 0 | 1 | 0 | 0 | 1 |
| 7  | 0 | 1 | 0 | 0 | 0 | 10 | 0 | 1 | 0 | 1 | 0 |
| 8  | 1 | 0 | 1 | 0 | 1 | 3  | 1 | 0 | 0 | 0 | 0 |
| 9  | 1 | 0 | 0 | 1 | 0 | 7  | 0 | 1 | 0 | 0 | 0 |
| 10 | 0 | 1 | 0 | 1 | 0 | 9  | 1 | 0 | 0 | 1 | 0 |

**Centers $a_k$ and Degree of heterogeneity $\varepsilon_k$**

|       | a | b | c | d | e |                  | a    | b    | c | d    | e    |
|-------|---|---|---|---|---|------------------|------|------|---|------|------|
| $a_1$ | 1 | 0 | 1 | 0 | 1 | $\varepsilon_1$  | 0    | 0    | 0 | 0    | 0.33 |
| $a_2$ | 0 | 1 | 0 | 1 | 0 | $\varepsilon_2$  | 0    | 0    | 0 | 0.25 | 0.5  |
| $a_3$ | 1 | 0 | 0 | 0 | 0 | $\varepsilon_3$  | 0.33 | 0.33 | 0 | 0.33 | 0    |

### CEM for the simplest model $[\varepsilon]$

- Exercise: When the proportions are supposed equal the classification log-likelihood to maximize

$$L_C(\mathbf{z}; \Theta) = L(\Theta; \mathbf{X}, \mathbf{z}) = \log\left(\frac{\varepsilon}{1-\varepsilon}\right) \sum_{i=1}^{n} \sum_{k=1}^{g} z_{ik} D(\mathbf{x}_i, \mathbf{a}_k) + np \log(1-\varepsilon)$$

  where $D(\mathbf{x}_i, \mathbf{a}_k) = \sum_{j=1}^{p} |x_{ij} - a_{kj}|$

- The parameter $\varepsilon$ is fixed for each cluster and for each variable, as $(\log(\frac{\varepsilon}{1-\varepsilon}) \leq 0)$ this maximization leads to the minimization of

$$W(\mathbf{z}, \mathbf{A}) = \sum_{i=1}^{n} \sum_{k=1}^{g} z_{ik} D(\mathbf{x}_i, \mathbf{a}_k)$$

- Exercise: The CEM algorithm is equivalent to the dynamical clustering method

### CEM and EM for the other models

- Exercise: Describe the different steps of CEM for the models $[\varepsilon_j]$, $[\varepsilon_k]$ and $[\varepsilon_{kj}]$
- Exercise: Deduce the different steps of EM for these models

**Nominal categorical data**

- Categorical data are a generalization of binary data
- Generally this kind of data is represented by a *complete disjunctive table* where the categories are represented by their indicators
- A variable $j$ with $h$ categories is represented by a binary vector such as

$$\left\{ \begin{array}{ll} x_i^{jh} = 1 & \text{if } i \text{ takes the categorie } h \text{ for} j \\ x_i^{jh} = 0 & \text{otherwise} \end{array} \right.$$

- The probability of the mixture can be written

$$f(\mathbf{x}_i; \Theta) = \sum_k \pi_k \prod_{j,h} (\alpha_k^{jh})^{x_{ij}}$$

where $\alpha_k^{jh}$ is the probability that the variable $j$ takes the categorie $h$ when an object belongs to the cluster $k$.

**Notation**

- $d_k^{jh} = \sum_i z_{ik} x_i^{jh}$
- $d^{jh} = \sum_i x_i^{jh}$
- $d_k = \sum_{j,h} d_k^{jh}$
- $d = \sum_k d_k = \sum_{k,j,h} x_i^{jh} = np$

**Example**

|    | a | b |
|----|---|---|
| 1  | 1 | 2 |
| 2  | 3 | 2 |
| 3  | 2 | 3 |
| 4  | 1 | 1 |
| 5  | 1 | 2 |
| 6  | 3 | 2 |
| 7  | 3 | 3 |
| 8  | 1 | 1 |
| 9  | 2 | 2 |
| 10 | 2 | 3 |

|    | a1 | a2 | a3 | b1 | b2 | b3 |
|----|----|----|----|----|----|----|
| 1  | 1  | 0  | 0  | 0  | 1  | 0  |
| 2  | 0  | 0  | 1  | 0  | 1  | 0  |
| 3  | 0  | 1  | 0  | 0  | 0  | 1  |
| 4  | 1  | 0  | 0  | 1  | 0  | 0  |
| 5  | 1  | 0  | 0  | 0  | 1  | 0  |
| 6  | 0  | 0  | 1  | 0  | 1  | 0  |
| 7  | 0  | 0  | 1  | 0  | 0  | 1  |
| 8  | 1  | 0  | 0  | 1  | 0  | 0  |
| 9  | 0  | 1  | 0  | 0  | 1  | 0  |
| 10 | 0  | 1  | 0  | 0  | 0  | 1  |

|    | a1 | a2 | a3 | b1 | b2 | b3 |
|----|----|----|----|----|----|----|
| 3  | 0  | 1  | 0  | 0  | 0  | 1  |
| 7  | 0  | 0  | 1  | 0  | 0  | 1  |
| 9  | 0  | 1  | 0  | 0  | 1  | 0  |
| 10 | 0  | 1  | 0  | 0  | 0  | 1  |
| 1  | 1  | 0  | 0  | 0  | 1  | 0  |
| 4  | 1  | 0  | 0  | 1  | 0  | 0  |
| 5  | 1  | 0  | 0  | 0  | 1  | 0  |
| 8  | 1  | 0  | 0  | 1  | 0  | 0  |
| 2  | 0  | 0  | 1  | 0  | 1  | 0  |
| 6  | 0  | 0  | 1  | 0  | 1  | 0  |

- $d_1^{a1} = 0, d_1^{a2} = 3,\ d_1^{a3} = 1,\ d_1^{b1} = 0, d_1^{b2} = 1,\ d_1^{b3} = 3$
- $d_1 = 8,\ d_2 = 8,\ d_3 = 4$
- $d = 8 + 8 + 4 = 10 \times 2$

**Interpretation of the model**

- The different steps of EM algorithm

  1. E-step: compute $s_{ik}$
  2. M-step: $\alpha_k^{jh} = \frac{\sum_i s_{ik} x_i^{jh}}{\sum_i s_{ik}}$ and $\pi_k = \frac{\sum_{i,k} s_{ik}}{n}$

- The different steps of CEM algorithm

  1. E-step: compute $s_{ik}$
  2. C-step: compute $\mathbf{z}$
  3. M-step (Exercise) : $\alpha_k^{jh} = \frac{\sum_i z_{ik} x_i^{jh}}{\sum_i z_{ik}} = \frac{d_k^{jh}}{\#z_k}$ and $\pi_k = \frac{\#z_k}{n}$

**Interpretation of the model**

- The classification log-likelihood can be written as

$$L_C(\mathbf{z}; \Theta) = \sum_{k,j,h} d_k^{jh} \log(\alpha_k^{jh}) + \sum_k \#z_k \log(\pi_k)$$

- When the proportions are supposed equal, the restricted likelihood

$$L_{CR}(\mathbf{z}; \Theta) = \sum_{k,j,h} d_k^{jh} \log(\alpha_k^{jh})$$

Given $\alpha_k^{jh} = \frac{d_k^{jh}}{\#z_k}$, it can be shown that CEM maximizes the mutual information

$$I(\mathbf{z}, J) = \sum_{k,j,h} \frac{d_k^{jh}}{d} \log \frac{d_k^{jh} d}{d_k d^{jh}}$$

This expression is very close to

$$\chi^2(\mathbf{z}, J) = \sum_{k,j,h} \frac{(d_k^{jh} d - d_k d^{jh})^2}{d_k d^{jh} d}$$

- Assuming that **X** derives form the latent class model whith equal proportions, the maximization of $L_C(\mathbf{z}; \Theta)$ is approximatively equivalent to use k-means with the $\chi^2$

**Parsimonious model**

- Number of the parameters in latent class model is equal $(g - 1) + g \times \sum_j (m^j - 1)$ where $m_i$ is the number of categories of $j$
- This number is smaller than $\prod_j m^j$ required by the complete log-linear model, example ($p = 10$, $g=5$, $m^j = 4$ for each $j$), this number is equal to $(5 - 1) + 5 * (40 - 10) = 154$
- This number can reduced by using parsimonious model by imposing constraints on the paremetre $\alpha_{kj}$. Instead to have a probability for each categorie, we associate for a categorie of $j$ having the same of value that the center for $j$ the probability $(1 - \varepsilon_{kj})$ and the other categories the probability $\varepsilon_{kj}/(m^j - 1)$
- Then the distribution depends on $\mathbf{a}_k$ and $\varepsilon_k$ defined by

$$\begin{cases} (1 - \varepsilon_{kj}) & \text{for } x_i^j = a_k^j \\ \varepsilon_{kj}/(m^j - 1) & \text{for } x_i^j \neq a_k^j \end{cases}$$

- The parametrization concerns only the variables instead of all categories
- This model is an extension of the Bernoulli model

## Example

```
library(Rmixmod)
data(birds)
dim(birds)
birds
xem.birds <- mixmodCluster(birds, 2)
summary(xem.birds)
****************************************
* number of modalities = 2 4 5 5 3
*** Cluster 1
* proportion = 0.6544
* center = 1.0000 3.0000 1.0000 1.0000 1.0000
* scatter = | 0.4937 0.4937 |
| 0.0761 0.0063 0.1741 0.0917 |
| 0.1521 0.1391 0.0043 0.0043 0.0043 |
| 0.0390 0.0045 0.0043 0.0259 0.0043 |
| 0.0577 0.0288 0.0289 |
*** Cluster 2
* proportion = 0.3456
* center = 2.0000 2.0000 2.0000 2.0000 1.0000
* scatter = | 0.4280 0.4280 |
| 0.1203 0.1463 0.0153 0.0107 |
| 0.0509 0.0751 0.0080 0.0080 0.0080 |
| 0.3641 0.5495 0.1288 0.0485 0.0080 |
| 0.1074 0.0940 0.0134 |
****************************************
```

**The simplest model**

- We assume that $(1 - \varepsilon_{kj})$ does not depend the cluster $k$ and the variable $j$

$$\begin{cases} (1 - \varepsilon) & \text{for } x_i^j = a_k^j \\ \varepsilon/(m^j - 1) & \text{for } x_i^j \neq a_k^j \end{cases}$$

- The restricted classification log-likelihood takes the following form

$$L_{CR}(\mathbf{z}; \Theta) = L(\Theta; \mathbf{X}, \mathbf{Z}) = \sum_{i,k} z_{ik} \left( \sum_j \log(\frac{\varepsilon}{1-\varepsilon}(m^j - 1))\delta(\mathbf{x}_i, \mathbf{a}_k) \right) + np \log(1 - \varepsilon)$$

or,

$$L_{CR}(\mathbf{z}; \Theta) = \sum_k \sum_{i \in z_k} d(\mathbf{x}_i, \mathbf{a}_k) + np \log(1 - \varepsilon)$$

where $D(\mathbf{x}_i, \mathbf{a}_k) = \sum_j \log(\frac{1-\varepsilon}{\varepsilon}(m^j - 1))\delta(x_{ij}, a_{kj})$

- If all variables have the same number of categories, the criterion to minimize is $\sum_{i,k} z_{ik} D(\mathbf{x}_i, \mathbf{a}_k)$, why ?
- The CEM is an extension of $k$-modes

**Contingency table**

- We can associate a multinomial model (Govaert and Nadif 2007), then the density of the model $\varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k) = B \sum_k \pi_k \alpha_{k1}^{x_{i1}} \ldots \alpha_{kp}^{x_{ip}}$ ( B does not depend on $\Theta$)

- Without $\log(B)$ we have $L_C(\mathbf{z}, \Theta) = \sum_i \sum_k z_{ik} \left( \log \pi_k + \sum_j x_{ij} \log(\alpha_{kj}) \right)$

- The mutual information quantifying the information shared between $\mathbf{z}$ and $J$:

$$I(\mathbf{z}, J) = \sum_{k,j} f_{kj} \log(\frac{f_{kj}}{f_k. f_{.j}})$$

- We have the relation $\sum_{k,j} \frac{(f_{kj} - f_k. f_{.j})^2}{f_k. f_{.j}} = \sum_{k,j} \frac{(f_{kj})^2}{f_k. f_{.j}} - 1$

- Using the following approximation : $x^2 - 1 \approx 2x \log(x)$ excellent in the neighborhood of 1 and good in $[0, 3]$, we have

$$\sum_{k,j} \frac{(f_{kj})^2}{f_k. f_{.j}} - 1 = \sum_{k,j} f_k. f_{.j} \left( (\frac{f_{kj}}{f_k. f_{.j}})^2 - 1 \right) \approx 2 \sum_{k,j} f_{kj} \log(\frac{f_{kj}}{f_k. f_{.j}})$$

- Th $I(\mathbf{z}, J) \approx \frac{1}{2N} \chi^2(\mathbf{z}, J)$

- When theis leads to proportions are assumed equal, the maximization of $L_C(\mathbf{z}, \Theta)$ is equivalent to the maximization of $I(\mathbf{z}, J)$ and approximately equivalent to the maximization of $\chi^2(\mathbf{z}, J)$

**The Gaussian model**

- The density can be written as: $f(\mathbf{x}_i; \Theta) = \sum_k \pi_k \varphi(\mathbf{x}_i; \boldsymbol{\mu}_k, \Sigma_k)$ where

$$\varphi(\mathbf{x}_i; \boldsymbol{\mu}_k, \Sigma_k) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_k|^{\frac{1}{2}}} exp\{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x}_i - \mu_k)\}$$

- Spectral decomposition of the variance matrix

$$\Sigma_k = \lambda_k D_k A_k D_k^T$$

  - $\lambda_k = |\Sigma_k|^{1/p}$ positive real represents the volume of the $k$th component
  - $A_k = Diag(a_{k1}, \ldots, a_{kp})$ formed by the normalized eigenvalues in decreasing order $|A_k| = 1$. It defines the shape of the $k$th cluster
  - $D_k$ formed by the eigenvectors. It defines the direction of the $k$th cluster

**Different Gaussian models**

- The Gaussian mixture depends on: proportions, centers, volumes, shapes and Directions then different models can be proposed
- In the following models proportions can be assumed equal or not

  1. Spherical models: $A_k = I$ then $\Sigma_k = \lambda_k I$. Two models $[\lambda I]$ and $[\lambda_k I]$
  2. Diagonal models: no constraint on $A_k$ but $D_k$ is a permutation matrix with $B_k = D_k A_k D_k^T$ such as $|B_k| = 1$, $\Sigma_k$ is diagonal. Four models $[\lambda B]$, $[\lambda_k B]$, $[\lambda B_k]$ and $[\lambda_k B_k]$
  3. General models: the eight models assuming equal or not volumes, shapes and directions $[\lambda DAD^T]$, $[\lambda_k DAD^T]$, $[\lambda DA_k D^T]$, $[\lambda_k DA_k D^T]$, $[\lambda D_k AD_k^T]$, $[\lambda_k D_k AD_k^T]$, $[\lambda D_k A_k D_k^T]$ and $[\lambda_k D_k A_k D_k^T]$

- Finally we have 28 models, we will study the problem of the choice of the models
- See for instance **mclust** and **Rmixmod**.

  mclust 5: Clustering, Classification and Density Estimation Using Gaussian Finite Mixture Models. by Luca Scrucca, Michael Fop, T. Brendan Murphy and Adrian E. Raftery, 2016

## CEM

- In clustering step, each $x_i$ is assigned to the cluster maximizing $s_{ik} \propto \pi_k \varphi(x_i; \mu_k, \Sigma_k)$ or equivalently the cluster that minimizes

$$-\log(\pi_k \varphi(x_i; \alpha_k)) = (x_i - \mu_k)^T \Sigma_k^{-1}(x_i - \mu_k) + \log|\Sigma_k| - 2\log(\pi_k) + cste$$

- From density to Distance (or dissimilarity), $x_i$ is assigned to the cluster according the following dissimilarity

$$D^2_{\Sigma_k^{-1}}(x_i; \mu_k) + \log|\Sigma_k| - 2\log(\pi_k)$$

where $D^2_{\Sigma_k^{-1}}(x_i; \mu_k) = (x_i - \mu_k)^T \Sigma_k^{-1}(x_i - \mu_k)$ is the Mahanalobis distance

- Note that when the proportions are supposed equal and the variances identical, the assignation is based only on

$$D^2_{\Sigma_k^{-1}}(x_i; \mu_k)$$

- When the proportions are supposed equal and for the spherical model $[\lambda I]$ ($\Sigma_k = I$), one uses the usual euclidean distance

$$D^2(x_i; \mu_k)$$

# library {mclust}

### Example

```
library(mclust)
data(diabetes)
class <- diabetes$class
table(class)
# class
# Chemical Normal Overt
# 36 76 33
X <- diabetes[,-1]
head(X)
res.pca=PCA(X)
clPairs(X, class)
res.mclust <- Mclust(X,3)
summary(res.mclust)
table(res.mclust$class,diabetes$class)
res.kmeans=kmeans(X,3,nstart=100)
table(res.kmeans$cluster,diabetes$class)
```

**Description of CEM**

- E-step: classical, C-step: Each cluster $z_k$ i formed by using $D^2(x_i; \mu_k)$
- M-step: Given the partition $\mathbf{z}$, we have to determine the parameter $\theta$ maximizing

$$L_C(\mathbf{z}, \Theta) = L(\Theta; \mathbf{X}, \mathbf{z}) = \sum_{i,k} z_{ik} \log\left(\pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k)\right) = \sum_k \sum_{i \in z_k} \log\left(\pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k)\right)$$

For the Gaussian model

$$-\frac{1}{2} \sum_k \left( \sum_i z_{ik}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k) + \#z_k \log|\Sigma_k| - 2\#z_k \log(\pi_k) \right)$$

- The parameter $\boldsymbol{\mu}_k$ is thus necessary the center $\boldsymbol{\mu}_k = \frac{\sum_i z_{ik} \mathbf{x}_i}{\#z_k}$
- The proportions satisfy $\pi_k = \frac{\#z_k}{n}$
- The parameters must then for the general model

$$F(\Sigma_1, \ldots, \Sigma_K) = \sum_k (\text{trace}(W_k \Sigma_k^{-1}) + \#z_k \log|\Sigma_k|)$$

where $W_k = \sum_i z_{ik}(\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$

**Consequence for the spherical model $[\lambda I]$**

- The function to maximize for the model $[\lambda I]$ becomes

$$F(\lambda) = \frac{1}{\lambda}\text{trace}(W) + np\log(\lambda)$$

where $W = \sum_k W_k$

With $\lambda = \frac{\text{trace}(W)}{np}$ maximizing $F(\lambda)$, the classification log-likelihood becomes

$$L_C(\mathbf{z};\Theta) = -\frac{np}{2}\text{trace}(W) + cste = -\frac{np}{2}W(\mathbf{z}) + cste$$

- Maximizing $L_C$ is equivalent to minimize the SSQ criterion minimized by the $k$means algorithm
- Interpretation
    - The use of the model $[\lambda I]$ assumes that the clusters are spherical having the same proportion and the same volume
    - The CEM is therefore an extension of the $k$means

### Description of EM

- E-step: classical
- M-step: we have to determine the parameter $\Theta$ maximizing $Q(\Theta, \Theta')$ taking the following form

$$L_C(\mathbf{z}; \Theta) = L(\Theta; \mathbf{X}, \mathbf{z}) = \sum_{i,k} s_{ik} \log \left( \pi_k \varphi(\mathbf{x}_i; \boldsymbol{\alpha}_k) \right)$$

For the Gaussian model

$$-\frac{1}{2} \sum_{i,k} \left( s_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) + s_{ik} \log |\Sigma_k| - 2s_{ik} \log(\pi_k) \right)$$

- The parameter $\boldsymbol{\mu}_k$ is thus necessary the center $\boldsymbol{\mu}_k = \frac{\sum_i s_{ik} \mathbf{x}_i}{\sum_i s_{ik}}$
- The proportions satisfy $\pi_k = \frac{\sum_i s_{ik}}{n}$
- The parameters $\Sigma_k$ must then minimize

$$F(\Sigma_1, \ldots, \Sigma_K) = \sum_k (\text{trace}(W_k \Sigma_k^{-1}) + \#z_k \log |\Sigma_k|)$$

where $W_k = \sum_i s_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$

## Von-Mises Fisher Mixture model

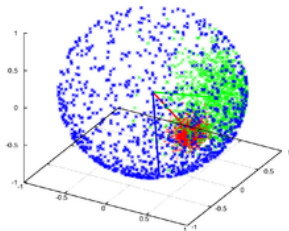### The von Mises-Fisher distribution (vMF)

Let $\mathbf{x}_i \in \mathbb{S}^{d-1}$ be a data point following a vMF distribution, then its pdf is

$$f(\mathbf{x}_i | \boldsymbol{\mu}, \kappa) = c_d(\kappa) \exp^{\kappa \boldsymbol{\mu}^\top \mathbf{x}_i}, \qquad (1)$$

$\boldsymbol{\mu}$: centroid parameter, $\kappa$: concentration parameter, such that $\|\boldsymbol{\mu}\| = 1$ and $\kappa \geq 0$. $c_d(\kappa) = \frac{\kappa^{\frac{d}{2}-1}}{(2\pi)^{\frac{d}{2}} I_{\frac{d}{2}-1}(\kappa)}$ $I_r(\kappa)$: the modified

Bessel function of the first kind and order $r$.



**Figure:** Impact of $\kappa$. blue: $\kappa = 1$, green: $\kappa = 10$, red: $\kappa = 100$

### The Mixture of vMF distributions (movMFs)

The data points $\mathbf{x}_1, \ldots, \mathbf{x}_n$ are supposed to be i.i.d and generated from a mixture of $g$ vMF distributions, with pdf:

$$f(\mathbf{x}_i | \Theta) = \sum_k \pi_k f_k(\mathbf{x}_i | \boldsymbol{\mu}_k, \kappa_k), \qquad (2)$$

where $\Theta = \{\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_g, \alpha_1, \ldots, \alpha_g, \kappa_1, \ldots, \kappa_g\}$

## Algorithms

**Log-likelihood**

$$L(\Theta; \mathbf{X}) = \sum_i \log \left( \sum_k \pi_k f(\mathbf{x}_i | \boldsymbol{\mu}_k, \kappa_k) \right),$$

**Complete data log-likelihood**

$$
\begin{aligned}
L_C(\mathbf{z}; \Theta) &= \sum_{i,k} z_{ik} \log \pi_k + \sum_{i,k} z_{ik} \log c_d(\kappa_k) + \sum_{i,k} z_{ik} \kappa_k \boldsymbol{\mu}_k^\top \mathbf{x}_i \\
&= \sum_{i,k} z_{ik} \log \pi_k + \sum_{i,k} z_{ik} \log c_d(\pi_k) + \sum_{i,k} z_{ik} \kappa_k \cos(\boldsymbol{\mu}_k, \mathbf{x}_i)
\end{aligned}
$$

**EM**

- E-step: finds the conditional expectation $\tilde{z}_{ik} = \mathbb{E}(z_{ik} = 1 | \mathbf{x}_i, \Theta^{(t)})$
- M-step: finds the new parameters $\Theta^{(t+1)}$ maximizing
  $Q(\Theta, \Theta^{(t)}) = \mathbb{E}\left( L(\Theta; \mathbf{X}, \mathbf{z}) | \mathbf{X}, \Theta^{(t)} \right)$ s.t. $\sum_k \pi_k = 1$, $\|\boldsymbol{\mu}_k\| = 1$ and $\kappa_k > 0$

**Hypotheses:** $\forall k, \pi_k = 1/g$ and $\kappa_k = \kappa$ the maximization of $L_C(\mathbf{z}; \Theta)$ and $\sum_{i,k} z_{ik} \cos(\mathbf{x}_i, \boldsymbol{\mu}_k)$ are equivalent

## Stochastic EM "SEM", (Celeux and Diebolt, 1985)

### Steps of SEM

- S-step between E-step and M-step
- In CEM (C-step), In SEM (S-step)
    - E-step: compute the posterior probabilities
    - S-step: This stochastic step consists to look for the partition $\bar{\mathbf{z}}$. Each object $i$ is assigned to the $k$th component. the parameter $k$ is selected according to the multinomial distribution $(s_{i1}, \ldots, s_{iK})$
    - M-step As the CEM algorithm this step is based on $\bar{\mathbf{z}}$

### Advantages and Disadvantages of SEM

- It gives good results when the size of data is large enough
- It can be used even if the number of clusters is unknown. It suffices to fix $K$ to $K_{max}$ the maximum number of clusters and this number can be reduced when the a cluster has a number of objects so lower that the estimation of parameters is not possible. For example when the cardinality of a cluster is less than a threshold, we run SEM with $(K - 1)$
- It can avoid the problem of initialization and other problems of EM
- Instability of the results. Solution: SEM (for estimation of paremetrs and the number of clusters), The obtained results are used by EM

# Stochastic Annealing EM "SAEM" (Celeux and Diebolt, 1992)

### Steps of SEM

- The aim of the SAEM is to reduce the "part" of random in estimations of the parameters
- SAEM is based on SEM and EM
- Solution
    - E-step: like for EM, SEM and CEM
    - S-step: like for SEM
    - M-step: The compute of parameters depends on this expression:

$$\boldsymbol{\theta}^{(t+1)} = \gamma^{(t+1)}\boldsymbol{\theta}_{SEM}^{(t+1)} + (1 - \gamma^{(t+1)})\boldsymbol{\theta}_{EM}^{(t+1)}$$

    The initial value of $\gamma = 1$ and decreases until 0.

## Outline

**Different approaches**

- In Finite mixture model, the problem of the choice of the model include the problem of the number of clusters

- To simplify the problem, we distinguish the two problems and we consider the model fixed and $K$ is unknown. Let be tow models $M_A$ and $M_B$. $\Theta(M_A)$ and $\Theta(M_B)$ indicates the "domain" of free parameters. if $L_{max}(M) = L(\hat{\theta}_M)$ where $\hat{\theta}_M = \text{argmax}\, L(\theta)$ then we have

$$\Theta(M_A) \subset \Theta(M_B) \Rightarrow L_{max}(M_A) \leq L_{max}(M_B)$$

For example $L_{max}[\pi_k \lambda_k I]_{K=2} \leq L_{max}[\pi_k \lambda_k I]_{K=3}$. Generally the likelihood increases with the number of clusters.

- First solution: Plot (Likelihood*number of clusters) and use the elbows

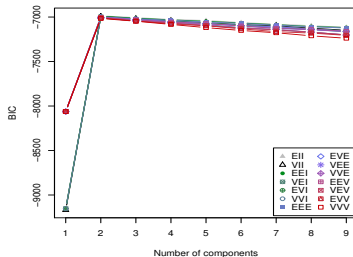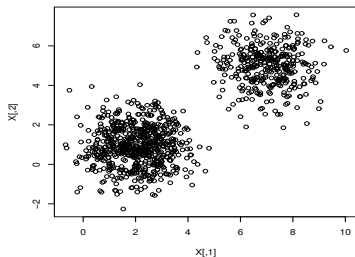- Second solution: Minimize the classical criteria (Criteria in competition) taking this form

$$C(M) = -2L_{max}(M) + \tau_C n_p(M)$$

where $n_p$ indicates the number of parameters of the model $M$, it represents the complexity of the model

- Different variants of this criterion AIC with $\tau_{AIC} = 2$, AIC3 with $\tau_{AIC} = 3$ and the famous

$$BIC(M) = -2L_{max}(M) + \log(n) n_p(M)$$

library(mclust)
res=Mclust(X)
plot(res)
summary(res)

# Outline

### Conclusion

- Finite mixture approach is interesting
- The CML approach gives interesting criteria and generalizes the classical criteria
- The different variants of EM offer good solutions
- The CEM algorithm is an extension of $k$-means and other variants
- The choice of the model is performed by using the maximum likelihood penalized by the number of parameters
- See **mclust** and **Rmixmod**
- Other Mixture models adapted to the nature of data (Text mining)