Machine Learning: Theory and Applications



Arnak S. Dalalyan American University of Armenia, Yerevan April 6, 2016



▲口→▲母→▲国→▲国→ 国 のへの



arnak-dalalyan.fr

1979 Born in Yerevan (tramvi park)



- 1979 Born in Yerevan (tramvi park)
- 1997 Bachelor degree in Mathematics, YSU
- 1999 Master degree in Probability and Statistics, Paris 6 University



- 1979 Born in Yerevan (tramvi park)
- 1997 Bachelor degree in Mathematics, YSU
- 1999 Master degree in Probability and Statistics, Paris 6 University
- 2001 PhD in Statistics, University of Le Mans (France)



- 1979 Born in Yerevan (tramvi park)
- 1997 Bachelor degree in Mathematics, YSU
- 1999 Master degree in Probability and Statistics, Paris 6 University
- 2001 PhD in Statistics, University of Le Mans (France)
- 2007 Habilitation in Statistics and Machine Learning, Paris 6 University



- 1979 Born in Yerevan (tramvi park)
- 1997 Bachelor degree in Mathematics, YSU
- 1999 Master degree in Probability and Statistics, Paris 6 University
- 2001 PhD in Statistics, University of Le Mans (France)
- 2007 Habilitation in Statistics and Machine Learning, Paris 6 University
- Now
- Professor of Statistics and Applied Mathematics, ENSAE ParisTech
 - Head of the Master in Data Science
 - Deputy chair of the Center for Data Science Paris-Saclay
 - Programme committee of the COLT (conf. on Learning Theory) and NIPS (Neural Information and Processing Systems), Associate Editor of EJS, JSPI, SISP and JSS.



What is Machine Learning?

The emphasis of machine learning is on devising automatic methods that perform a given task based on what was experienced in the past.

Here are several examples :

- optical character recognition : categorize images of handwritten characters by the letters represented
- topic spotting : categorize news articles (say) as to whether they are about politics, sports, entertainment, etc.
- medical diagnosis : diagnose a patient as a sufferer or non-sufferer of some disease
- customer segmentation : dividing a customer base into groups of individuals that are similar in specific ways relevant to marketing (such as gender, interests, spending habits)
- fraud detection : identify credit card transactions (for instance) which may be fraudulent in nature



Typical tasks of machine learning

Supervised learning :

- Prediction : use historical data for predicting future values
- Classification : identifying to which of a set of classes a new observation belongs, on the basis of a training set of data containing observations whose class membership is known
- Ranking : construction of ranking models for information retrieval systems

Unsupervised learning :

- Outlier detection : detecting a few observations that deviate so much from other observations as to arouse suspicion that it was generated by a different mechanism
- Clustering : grouping a set of objects into homogeneous groups based on some similarity measure
- Dimensionality reduction : mapping of data to a lower dimensional space such that uninformative variance in the



👾 🤅 🤄 a is discarded



Outline

Introduction to clustering

- K-means and partitioning around medoids
 - Theoretical background
 - R code
 - An example

Gaussian mixtures and EM algorithm

- Theoretical background
- R code
- An example
- Hierarchical Clustering
 - Theoretical background
 - R code
 - An example



Main objective

The goal of clustering analysis is to find high-quality clusters such that the inter-cluster similarity is low and the intra-cluster similarity is high.





The general diagram





The general diagram





The general diagram



In general, there is no objective measure of quality for a clustering algorithm. The satisfaction of the end-user is perhaps the best manner of evaluation.



Notation

• We are given *n* examples represented as a *p*-vector with real-values entries :

$$\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathbb{R}^p, \qquad \mathbf{X}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix}$$

• We denote by $\|\mathbf{X}_i - \mathbf{X}_j\|$ the (Euclidean) distance between two examples :

$$\|\mathbf{X}_i - \mathbf{X}_j\|^2 = \sum_{\ell=1}^p (x_{i,\ell} - x_{j,\ell})^2$$

For a set G, subset of {1,...,n}, we denote by X_G the average of the examples X_i corresponding to i ∈ G, that is

$$\overline{\mathbf{X}}_G = \frac{1}{|G|} \sum_{i \in G} \mathbf{X}_i.$$



Definition of the method

- Main idea : A good clustering algorithm divides the sample into *K* groups such that the variance within each group is mall.
- Mathematically speaking, this corresponds to solving with respect to $G_1, \ldots, G_K \subset \{1, \ldots, n\}$ and $\mathbf{C}_1, \ldots, \mathbf{C}_K \in \mathbb{R}^p$ the following optimisation problem :

$$\min_{G_1,\ldots,G_K} \min_{\mathbf{C}_1,\ldots,\mathbf{C}_K} \underbrace{\sum_{k=1}^K \sum_{i \in G_k} \|\mathbf{X}_i - \mathbf{C}_k\|^2}_{\Psi(G_{1:K},\mathbf{C}_{1:K})},$$

where G_1, \ldots, G_K runs over all possible partitions of $\{1, \ldots, n\}$.

• Important remark : the minimization should be done for a fixed value of *K* (prescribed number of clusters), otherwise the solution is obvious (and meaningless), isn't it?



How to minimize the cost?

• Cost function : $\Psi(G_{1:K}, \mathbf{C}_{1:K}) = \sum_{k=1}^{K} \sum_{i \in G_k} \|\mathbf{X}_i - \mathbf{C}_k\|^2$ is hard to minimize. It is a combinatorial optimization problem which is known to be NP-hard.



How to minimize the cost?

- Cost function : $\Psi(G_{1:K}, \mathbf{C}_{1:K}) = \sum_{k=1}^{K} \sum_{i \in G_k} \|\mathbf{X}_i \mathbf{C}_k\|^2$ is hard to minimize. It is a combinatorial optimization problem which is known to be NP-hard.
- Local minimum : finding a local minimum is easy.



How to minimize the cost?

- Cost function : $\Psi(G_{1:K}, \mathbf{C}_{1:K}) = \sum_{k=1}^{K} \sum_{i \in G_k} \|\mathbf{X}_i \mathbf{C}_k\|^2$ is hard to minimize. It is a combinatorial optimization problem which is known to be NP-hard.
- Local minimum : finding a local minimum is easy.
 - Step 1 For fixed centers $C_{1:K}$ the minimizer of $\Psi(G_{1:K}, C_{1:K})$ w.r.t. $G_{1:K}$ can be easily computed :



How to minimize the cost?

- Cost function : $\Psi(G_{1:K}, \mathbf{C}_{1:K}) = \sum_{k=1}^{K} \sum_{i \in G_k} \|\mathbf{X}_i \mathbf{C}_k\|^2$ is hard to minimize. It is a combinatorial optimization problem which is known to be NP-hard.
- Local minimum : finding a local minimum is easy.
 - Step 1 For fixed centers $C_{1:K}$ the minimizer of $\Psi(G_{1:K}, C_{1:K})$ w.r.t. $G_{1:K}$ can be easily computed :

► G_k contains all the points \mathbf{X}_i for which the closest point in the set { $\mathbf{C}_1, \ldots, \mathbf{C}_k$ } is \mathbf{C}_k .



How to minimize the cost?

- Cost function : $\Psi(G_{1:K}, \mathbf{C}_{1:K}) = \sum_{k=1}^{K} \sum_{i \in G_k} \|\mathbf{X}_i \mathbf{C}_k\|^2$ is hard to minimize. It is a combinatorial optimization problem which is known to be NP-hard.
- Local minimum : finding a local minimum is easy.
 - Step 1 For fixed centers $C_{1:K}$ the minimizer of $\Psi(G_{1:K}, C_{1:K})$ w.r.t. $G_{1:K}$ can be easily computed :

► G_k contains all the points \mathbf{X}_i for which the closest point in the set { $\mathbf{C}_1, \ldots, \mathbf{C}_k$ } is \mathbf{C}_k .

Step 2 For fixed groups $G_{1:K}$, the minimizer of $\Psi(G_{1:K}, \mathbf{C}_{1:K})$ w.r.t. $\mathbf{C}_{1:K}$ can be easily computed :



How to minimize the cost?

- Cost function : $\Psi(G_{1:K}, \mathbf{C}_{1:K}) = \sum_{k=1}^{K} \sum_{i \in G_k} \|\mathbf{X}_i \mathbf{C}_k\|^2$ is hard to minimize. It is a combinatorial optimization problem which is known to be NP-hard.
- Local minimum : finding a local minimum is easy.
 - Step 1 For fixed centers $C_{1:K}$ the minimizer of $\Psi(G_{1:K}, C_{1:K})$ w.r.t. $G_{1:K}$ can be easily computed :

► G_k contains all the points \mathbf{X}_i for which the closest point in the set { $\mathbf{C}_1, \ldots, \mathbf{C}_k$ } is \mathbf{C}_k .

Step 2 For fixed groups $G_{1:K}$, the minimizer of $\Psi(G_{1:K}, \mathbf{C}_{1:K})$ w.r.t. $\mathbf{C}_{1:K}$ can be easily computed :

$$\mathbf{V}_k = \overline{\mathbf{X}}_{G_k}.$$



How to minimize the cost?

- Cost function : $\Psi(G_{1:K}, \mathbf{C}_{1:K}) = \sum_{k=1}^{K} \sum_{i \in G_k} \|\mathbf{X}_i \mathbf{C}_k\|^2$ is hard to minimize. It is a combinatorial optimization problem which is known to be NP-hard.
- Local minimum : finding a local minimum is easy.
 - Step 1 For fixed centers $C_{1:K}$ the minimizer of $\Psi(G_{1:K}, C_{1:K})$ w.r.t. $G_{1:K}$ can be easily computed :

► G_k contains all the points \mathbf{X}_i for which the closest point in the set { $\mathbf{C}_1, \ldots, \mathbf{C}_k$ } is \mathbf{C}_k .

Step 2 For fixed groups $G_{1:K}$, the minimizer of $\Psi(G_{1:K}, \mathbf{C}_{1:K})$ w.r.t. $\mathbf{C}_{1:K}$ can be easily computed :

 $\mathbf{P} \mathbf{C}_k = \overline{\mathbf{X}}_{G_k}.$

• Iterative algorithm : initialize at some $C_{1:{\cal K}}$ and then alternate between the two aforementioned steps until the convergence.



Alternating minimization

An illustrative example :

Ideal clustering



Alternating minimization

An illustrative example :

K-means clustering



Initialization



Alternating minimization

An illustrative example :



interation 1 (a)



Alternating minimization

An illustrative example :



interation 1 (b)



Alternating minimization

An illustrative example :



interation 2 (a)



Alternating minimization

An illustrative example :



interation 2 (b)



Alternating minimization

An illustrative example :



interation 3 (a)



Alternating minimization

An illustrative example :



K-means clustering

interation 3 (a)

In general, a dozen iterations are sufficient for the algorithm to converge ...



Alternating minimization

An illustrative example :



K-means clustering

In general, a dozen iterations are sufficient for the algorithm to converge $\underline{to a \ local \ minimum.}$



Alternating minimization

Not all initializations are good :

 $\begin{array}{c} \mathbf{N} \\ \mathbf{$



Alternating minimization

Not all initializations are good :





Alternating minimization

Not all initializations are good :



interation 1 (a)



Alternating minimization

Not all initializations are good :



interation 1 (b)



Alternating minimization

Not all initializations are good :



interation 2 (a)



Alternating minimization

Not all initializations are good :

K-means clustering



Most implementations (and this is the case for the R function kmeans) compute the clusterings for several initializations and select the one having the minimal cost.



The iris dataset

The Iris flower data set was introduced by Ronald Fisher in his 1936 paper. It is sometimes called Anderson's Iris data set because Edgar Anderson collected the data to quantify the morphologic variation of Iris flowers of three related species.



Fisher's <i>Iris</i> Data					
Sepal length +	Sepal width +	Petal length +	Petal width +	Species +	
5.1	3.5	1.4	0.2	I. setosa	
4.9	3.0	1.4	0.2	I. setosa	



K-means clustering R code



Here is a toy example of clustering with **R**.





Component 1 These two components explain 95.81 % of the poin



K-means clustering R code



Here is a toy example of clustering with **R**.





Component 1 These two components explain 95.81 % of the poin



R code



Pay attention to :

- use several initializations for reducing the randomness of the result
- normalize the columns of the data matrix







Remarks

- Breaking the ties In the step of assigning data points to clusters it may happen that two centers are the distance from a data point. Such a tie is usually broken at random using a coin toss.
- K-medians If the data is likely to contain outliers, it might be better to replace the squared Euclidean distance by the "manhattan" distance :

$$\|\mathbf{X}_{i} - \mathbf{X}_{j}\|_{1} = \sum_{\ell=1}^{p} |x_{i\ell} - x_{j\ell}|$$

and to minimize the cost function

$$\sum_{k=1}^{K}\sum_{i\in G_k}\|\mathbf{X}_i-\mathbf{C}_k\|_1.$$

This method is referred to as k-medians since C_k is necessarily the median of the cluster $\{X_i : i \in G_k\}$.

 PAM More generally, one can consider any measure of dissimilarity between data points. The resulting algorithm is called partitioning around medoids.



Clustering CAC40 companies

Let us consider the problem of clustering the CAC40 companies according to their stock prices during the past 2 months.

- Each company is described by its daily stock returns (44 real values), the intraday variation (45 positive values) and the daily volume of exchanges (45 positive values). Data is downloaded from http://www.abcbourse.com/.
- We first determine the number of clusters using the function kmeans runs



Clustering CAC40 companies

The result of k-means

Cluster 1 Size :23	Cluster 2 Size :7	Cluster 3 Size :10
Total L'oreal Lvmh Schneider El Veolia Env GDF Suez Alstom EDF Pernod Ric Danone	Credit Agr Alcatel-Lucent Societe Generale Bnp Paribas Renault Orange Arcelor Mittal	Accor Bouygues Lafarge Michelin Saint Gobain Vinci Valeo Publicis Groupe Technip Gemalto



Expectation maximization for Gaussian mixtures

Gaussian Mixture (GM) Model based clustering Theoretical background

- ► The data points X₁,..., X_n are assumed to be generated at random according to the following mechanism :
 - *n* "lablels" Z₁,..., Z_n are independently generated according to a distribution π on the set {1,..., K} (if Z_i = k then X_i belongs to the kth cluster).
 - for each i = 1, ..., n, if $Z_i = k$, then the data point \mathbf{X}_i is drawn from a multivariate Gaussian distribution with a mean $\boldsymbol{\mu}_k$ and a covariance matrix $\boldsymbol{\Sigma}_k$.
- ▶ Only $X_1, ..., X_n$ are observed. The goal is to find labels $\hat{Z}_1, ..., \hat{Z}_n$ such that the probability of the event $\{\hat{Z}_i = Z_i, \forall i\}$ is as high as possible.



From a mathematical point of view, X_1, \ldots, X_n are independent and have the density :

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \varphi(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

with $\varphi(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma})$ being the Gaussian density

$$\varphi(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{p/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\Big\{-\frac{1}{2}\|\boldsymbol{\Sigma}^{-1/2}(\boldsymbol{x}-\boldsymbol{\mu})\|^2\Big\}.$$

Gaussian density in 2D Density of a GM in 2D (K = 2)



From clustering to estimation

Since X_1, \ldots, X_n are independent and have the density :

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} \pi_k \varphi(\boldsymbol{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

if the parameters $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ of the model were known, we would determine the cluster number assigned to **x** by maximizing w.r.t. k

$$\mathbf{P}(Z=k|\mathbf{X}=\mathbf{x})=\frac{\pi_k\varphi(\mathbf{x}|\boldsymbol{\mu}_k,\boldsymbol{\Sigma}_k)}{p(\mathbf{x})}.$$



From clustering to estimation

Since X_1, \ldots, X_n are independent and have the density :

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} \pi_k \varphi(\boldsymbol{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

if the parameters $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ of the model were known, we would determine the cluster number assigned to **x** by maximizing w.r.t. k

$$\mathbf{P}(Z=k|\mathbf{X}=\mathbf{x})=\frac{\pi_k\varphi(\mathbf{x}|\boldsymbol{\mu}_k,\boldsymbol{\Sigma}_k)}{p(\mathbf{x})}.$$

► Therefore, it is natural to first estimate $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ by $\{\hat{\pi}, (\hat{\mu}_k)_k, (\hat{\Sigma}_k)_k\}$ and then to set

$$\hat{Z} = rg\max_{k=1,...,K} \hat{\pi}_k arphi(oldsymbol{x} | \hat{oldsymbol{\mu}}_k, \hat{\Sigma}_k).$$



From clustering to estimation

Since X_1, \ldots, X_n are independent and have the density :

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} \pi_k \varphi(\boldsymbol{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

if the parameters $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ of the model were known, we would determine the cluster number assigned to **x** by maximizing w.r.t. k

$$\mathbf{P}(Z=k|\mathbf{X}=\mathbf{x})=\frac{\pi_k\varphi(\mathbf{x}|\boldsymbol{\mu}_k,\boldsymbol{\Sigma}_k)}{p(\mathbf{x})}.$$

► Therefore, it is natural to first estimate $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ by $\{\hat{\pi}, (\hat{\mu}_k)_k, (\hat{\Sigma}_k)_k\}$ and then to set

$$\hat{Z} = \arg \max_{k=1,...,K} \hat{\pi}_k \varphi(\mathbf{x} | \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}}_k).$$

Problem : how to estimate these parameters?

Estimation problem

Estimate the quantities $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ based on the observation of n independent random variables with density $p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \varphi(\mathbf{x} | \mu_k, \Sigma_k)$.

Main difficulty : the maximum likelihood estimator

$$\left\{\hat{\boldsymbol{\pi}}, \left(\hat{\boldsymbol{\mu}}_{k}\right)_{k}, \left(\hat{\boldsymbol{\Sigma}}_{k}\right)_{k}
ight\}^{ML} = \arg\max_{\left\{\boldsymbol{\pi}, \left(\boldsymbol{\mu}_{k}\right)_{k}, \left(\boldsymbol{\Sigma}_{k}\right)_{k}
ight\}} \sum_{i=1}^{n}\log p(\mathbf{X}_{i})$$

is, in practice, impossible to compute.



Estimation problem

Estimate the quantities $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ based on the observation of n independent random variables with density $p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \varphi(\mathbf{x} | \mu_k, \Sigma_k)$.

Main difficulty : the maximum likelihood estimator

$$\left\{\hat{\boldsymbol{\pi}}, (\hat{\boldsymbol{\mu}}_k)_k, (\hat{\boldsymbol{\Sigma}}_k)_k\right\}^{ML} = \arg\max_{\left\{\boldsymbol{\pi}, (\boldsymbol{\mu}_k)_k, (\boldsymbol{\Sigma}_k)_k\right\}} \sum_{i=1}^n \log\left\{\sum_{k=1}^K \pi_k \varphi(\mathbf{X}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right\}$$

is, in practice, impossible to compute.



EM algorithm : main idea

Estimation problem

Estimate the quantities $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ based on the observation of n independent random variables with density $p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \varphi(\mathbf{x} | \mu_k, \Sigma_k)$.

Main difficulty : the maximum likelihood estimator

$$\left\{\hat{\boldsymbol{\pi}}, (\hat{\boldsymbol{\mu}}_k)_k, (\hat{\boldsymbol{\Sigma}}_k)_k\right\}^{ML} = \arg\max_{\left\{\boldsymbol{\pi}, (\boldsymbol{\mu}_k)_k, (\boldsymbol{\Sigma}_k)_k\right\}} \sum_{i=1}^n \log\left\{\sum_{k=1}^K \pi_k \varphi(\mathbf{X}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right\}$$

is, in practice, impossible to compute.

Work-around : use the identity

$$\log\left\{\sum_{k=1}^{K}\pi_{k}a_{k}\right\} = \max_{\nu}\sum_{k=1}^{K}\left(\nu_{k}\log(\pi_{k}a_{k}) - \nu_{k}\log\nu_{k}\right).$$



EM algorithm : main idea

Estimation problem

Estimate the quantities $\{\pi, (\mu_k)_k, (\Sigma_k)_k\}$ based on the observation of n independent random variables with density $p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \varphi(\mathbf{x} | \mu_k, \Sigma_k)$.

Main difficulty : the maximum likelihood estimator

$$\left\{\hat{\boldsymbol{\pi}}, (\hat{\boldsymbol{\mu}}_k)_k, (\hat{\boldsymbol{\Sigma}}_k)_k\right\}^{ML} = \arg\max_{\left\{\boldsymbol{\pi}, (\boldsymbol{\mu}_k)_k, (\boldsymbol{\Sigma}_k)_k\right\}} \sum_{i=1}^n \log\left\{\sum_{k=1}^K \pi_k \varphi(\mathbf{X}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right\}$$

is, in practice, impossible to compute.

Work-around : use the identity

$$\log\Big\{\sum_{k=1}^{K}\pi_k a_k\Big\} = \max_{\substack{\nu \in \mathbb{R}_+^{K} \\ \sum \nu_k = 1}} \sum_{k=1}^{K} \bigg(\nu_k \log a_k - \nu_k \log(\nu_k/\pi_k)\bigg).$$



Let us denote $\Omega = \{\pi, (\mu_k), (\Sigma_k)\}$ and $\hat{\Omega}^{ML} = \{\hat{\pi}, (\hat{\mu}_k), (\hat{\Sigma}_k)\}^{ML}$. Using the previous formulae, we can write $\hat{\Omega}^{ML}$ as follows

$$\hat{\Omega}^{ML} = \arg\max_{\Omega} \max_{\nu_{i}} \underbrace{\sum_{i=1}^{n} \sum_{k=1}^{K} \left\{ \nu_{ik} \log(\pi_{k} \varphi(\mathbf{X}_{i} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})) - \nu_{ik} \log \nu_{ik} \right\}}_{\boldsymbol{G}(\Omega, \nu)}.$$

(1)



Let us denote $\Omega = \{\pi, (\mu_k), (\Sigma_k)\}$ and $\hat{\Omega}^{ML} = \{\hat{\pi}, (\hat{\mu}_k), (\hat{\Sigma}_k)\}^{ML}$. Using the previous formulae, we can write $\hat{\Omega}^{ML}$ as follows $\hat{\Omega}^{ML} = \arg\max_{\Omega} \max_{\nu} \boldsymbol{G}(\Omega, \nu).$ (1)



Let us denote $\Omega = \{\pi, (\mu_k), (\Sigma_k)\}$ and $\hat{\Omega}^{ML} = \{\hat{\pi}, (\hat{\mu}_k), (\hat{\Sigma}_k)\}^{ML}$. Using the previous formulae, we can write $\hat{\Omega}^{ML}$ as follows

$$\hat{\Omega}^{ML} = \arg\max_{\Omega} \max_{\boldsymbol{\nu}} \boldsymbol{\sigma}(\Omega, \boldsymbol{\nu}).$$
(1)

- For a fixed Ω , the maximum w.r.t. ν in (1) is explicitly computable.
- For a fixed ν, the maximum w.r.t. Ω in (1) is explicitly computable as well.



Let us denote $\Omega = \{\pi, (\mu_k), (\Sigma_k)\}$ and $\hat{\Omega}^{ML} = \{\hat{\pi}, (\hat{\mu}_k), (\hat{\Sigma}_k)\}^{ML}$. Using the previous formulae, we can write $\hat{\Omega}^{ML}$ as follows

$$\hat{\Omega}^{ML} = \arg\max_{\Omega} \max_{\boldsymbol{\nu}} \boldsymbol{\sigma}(\Omega, \boldsymbol{\nu}).$$
(1)

- For a fixed Ω , the maximum w.r.t. ν in (1) is explicitly computable.
- For a fixed ν, the maximum w.r.t. Ω in (1) is explicitly computable as well.

EM-algorithm

Initialize Ω . Then iteratively (until convergence)

- update $oldsymbol{
 u}$ by solving (1) with fixed Ω
- update Ω by solving (1) with fixed u.



EM algorithm : how it works

EM-algorithm

Initialize Ω . Then iteratively (until convergence)

- update $oldsymbol{
 u}$ by solving (1) with fixed Ω
- update Ω by solving (1) with fixed u.





EM algorithm : summary and remarks

- The goal of the EM-algorithm is to approximate the maximum likelihood estimator.
- The EM-algorithm solves a nonconvex optimization problem. Therefore, there is no guarantee that it finds the global optimum. Generally, it does not !
- ► To apply the EM algorithm, the sample size *n* should be significantly larger than $p \times k + p^2$.
- One can adapt the EM algorithm to other (non Gaussian) distributions.
- ► It is possible in the EM algorithm to assume that all the Σ_k 's are equal, or that they are all diagonal.
- ► The auxiliary values ν_{ik} computed by the EM-algorithm estimate the probability of **X**_i to belong to the cluster k.



The EM-algorithm is implemented in the R-package mclust.







Gaussian Mixture (GM) Model based clustering R code for CAC 40 dataset





Clustering CAC40 companies The result of EM

Cluster 2 Cluster 3 Cluster 4 Cluster 5 Cluster 1 Size :12 Size :10 Size :8 Size :8 Size :2 Safran Air Lig Credit Agr l'oreal Total Solvav Carrefour Michelin Accor Axa Sanofi Vivendi Bouyques Kering Danone Pernod Ric Alcatel-Luc. Unibail-Rod. Schneider El Lafarge Renault I vmh Soc. Gen. Valeo Veolia Env Essilor Intl Publicis Gr. Saint Gobain Bnp Paribas Cap Gemini Orange Technip Vinci Alstom GDF Suez Gemalto Airbus Group FDF Arcelor Mit. Legrand SA





~ ~ ~ ~

◆□▶◆□▶◆三▶◆三▶ ● 三臣