

Model-based clustering and classification

part 1: basics

Julien JACQUES

Université Lumière Lyon 2

Introduction

The mixture model

Mixture model estimation in classification

Mixture model estimation in clustering

Introduction

Clustering and classification

- ▶ **clustering** (unsupervised):

to group a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense or another) to each other than to those in other groups (clusters)

- ▶ **classification** (discrimination, scoring / supervised):

to predict the group of a new observation from a labeled sample

Clustering and classification

Notations

- ▶ observations are described by p features $\mathbf{X} = (X_1, \dots, X_p) \in E$
($E = \mathbb{R}^p, \dots$)
- ▶ $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})$ is the features for observation i ($1 \leq i \leq n$)
- ▶ $Z_i \in \{1, \dots, K\}$ is the group of observation i

Clustering versus classification

Clustering

- ▶ Z_i unknown
- ▶ goal: to predict Z_1, \dots, Z_n from $\mathbf{X}_1, \dots, \mathbf{X}_n$
- ▶ Z_1, \dots, Z_n are a posteriori interpreted in order to give significance to the clusters

Classification

- ▶ Z_i observed
- ▶ goal: to build a classification rule r from $(\mathbf{X}_1, Z_1), \dots, (\mathbf{X}_n, Z_n)$:

$$r : \mathbf{X} \longrightarrow r(\mathbf{X}) = Z$$

- ▶ to use this rule in order to classify new observation for which the group is unknown

Applications

Clustering

- ▶ exploratory analysis: to give a simplified representation of data in order to understand them
- ▶ example:
 - ▶ to recognize communities in social networks
 - ▶ to extract topics from corpus of documents
 - ▶ typology of customers in CRM (Customer Relationship Management)

Classification

- ▶ predictive analysis: to predict Z (categorical) from covariates \mathbf{X} (categorical, continuous...)
- ▶ example: to predict the probability (score) ...
 - ▶ marketing: ... for a new customer to buy a product
 - ▶ medicine: ... for a patient to be suffering from a disease
 - ▶ finance: ... for a firm to enter bankruptcy

Different methods

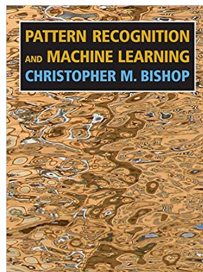
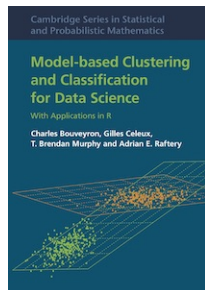
Clustering

- ▶ geometric methods
 - ▶ kmeans, hierarchical clustering
- ▶ probabilistic methods
 - ▶ **mixture models**

Classification

- ▶ generative methods: estimation of $p(\mathbf{X}, Z)$
 - ▶ **mixture models** (linear/quadratic discriminant analysis, ...)
- ▶ predictive methods: estimation of $p(Z|\mathbf{X})$
 - ▶ logistic regression, K- nearest neighbors, classification tree, SVM, neural networks. . .

References



- ▶ G. Celeux & G. Govaert (1995), *Gaussian parsimonious clustering models*, Pattern Recognition, 28(5), 781–793.
- ▶ L. Scrucca, M. Fop, T. B. Murphy & A.E. Raftery (2016), *mclust 5: Clustering, Classification and Density Estimation Using Gaussian Finite Mixture Models*, The R Journal, 8(1), 289–317.

The mixture model

Definition and notation

Idea: each group is described by its *own* probability distribution

$$\mathbf{X}|Z = k \sim f(\mathbf{x}, \theta_k) = f_k(\mathbf{x})$$

for instance, $f(\cdot, \theta_k)$ can be

- ▶ Continuous features: the Gaussian distrib. ($\theta_k = (\mu_k, \Sigma_k)$), the Student distribution for more heavy tails...
- ▶ Binary features: multivariate Bernoulli distrib., $\theta_k = (\alpha_{kj})_{1 \leq j \leq p}$
- ▶ Categorical features: multinomial distribution...

Definition and notation

mixing proportion

$$Z = k \Leftrightarrow \tilde{Z} = (0, \dots, 0, \underbrace{1}_{k\text{-th position}}, 0, \dots, 0)$$

$$\tilde{Z} \sim \mathcal{M}(1, p_1, \dots, p_K)$$

where $p_k = P(Z = k) = P(\tilde{Z}_k = 1)$ is the mixing proportion of group k

The mixture model

- ▶ marginal distribution of \mathbf{X} (mixture density)

$$\mathbf{X} \sim \sum_{k=1}^K p_k f_k(\mathbf{x}) = f_{\mathbf{X}}(\mathbf{x}).$$

The mixture model

- ▶ marginal distribution of \mathbf{X} (mixture density)

$$\mathbf{X} \sim \sum_{k=1}^K p_k f_k(\mathbf{x}) = f_{\mathbf{X}}(\mathbf{x}).$$

proof:

$$\begin{aligned} P(\mathbf{X} \in I) &= P(\mathbf{X} \in I \cap Z \in \{1, \dots, K\}) \\ &= \sum_{k=1}^K P(\mathbf{X} \in I \cap Z = k) \\ &= \sum_{k=1}^K P(\mathbf{X} \in I | Z = k) P(Z = k) \end{aligned}$$

The mixture model

- ▶ conditional probability that \mathbf{x} belongs to group k (via Bayes theorem):

$$P(Z = k | \mathbf{X} = \mathbf{x}) = \frac{p_k f_k(\mathbf{x})}{f_{\mathbf{X}}(\mathbf{x})} = t_k(\mathbf{x}).$$

The mixture model

- ▶ conditional probability that \mathbf{x} belongs to group k (via Bayes theorem):

$$P(Z = k | \mathbf{X} = \mathbf{x}) = \frac{p_k f_k(\mathbf{x})}{f_{\mathbf{X}}(\mathbf{x})} = t_k(\mathbf{x}).$$

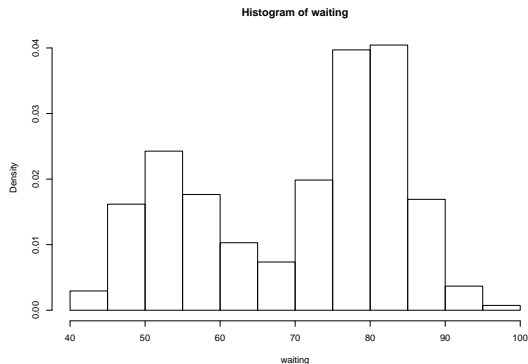
proof:

$$\begin{aligned} P(Z = k | \mathbf{X} = \mathbf{x}) &= \frac{f(\mathbf{x} | Z = k) P(Z = k)}{f(\mathbf{x})} \\ &= \frac{f(\mathbf{x} | Z = k) P(Z = k)}{\sum_{\ell=1}^K f(\mathbf{x} | Z = \ell) P(Z = \ell)} \end{aligned}$$

Example - Faithful

Waiting time between eruptions and the duration of the eruption for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA.

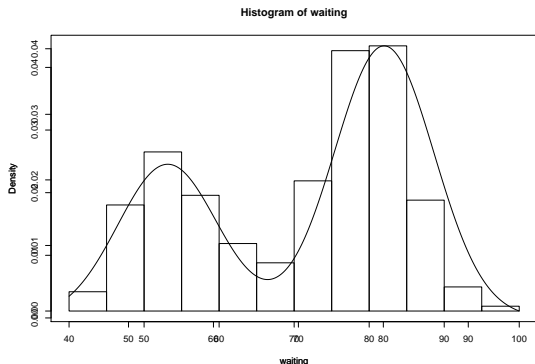
```
attach(faithful)
hist(waiting,prob=T)
```



Example - Faithful

Waiting time between eruptions and the duration of the eruption for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA.

```
library(mclust)
res=Mclust(waiting)
hist(waiting,prob=T);par(new=TRUE)
plot(res,'density')
```



Clustering and classification rule

Clustering and classification rule

Let assume that all the mixture model parameters (p_k and the parameters of f_k) are known (*they will be estimated in practice from data*)

Clustering: each observation \mathbf{x}_i is classified into the group k maximizing the conditional probability $t_k(\mathbf{x}_i) = P(Z = k | \mathbf{X} = \mathbf{x}_i)$:

$$Z = \operatorname{argmax}_k t_k(\mathbf{x}_i)$$

Classification: it will depend of the cost of misclassification (not necessary symmetric)

Classification rule

To define a classification rule

$$r : \mathbf{x} \in \mathbb{E} \rightarrow r(\mathbf{x}) \in \{1, \dots, K\}.$$

is equivalent to divide \mathbb{E} into K subsets Ω_k s.t.

$$\Omega_1 \cup \dots \cup \Omega_K = \mathbb{E},$$

$$\Omega_k \cap \Omega_\ell = \emptyset$$

and $\mathbf{x} \in \Omega_k \Leftrightarrow r(\mathbf{x}) = k.$

Probability of misclassification

Probability of classifying an observation of group G_k into G_ℓ ($\ell \neq k$) with r :

$$e_{k\ell}(r) = P(r(\mathbf{X}) = \ell | Z = k) = \int_{\Omega_\ell} f_k(\mathbf{x}) d\mathbf{x}.$$

Probability of misclassification of an observation of G_k with r :

$$e_k(r) = P(r(\mathbf{X}) \neq k | Z = k) = \sum_{\ell \neq k} e_{k\ell}(r) = \int_{\mathbb{E} \setminus \Omega_k} f_k(\mathbf{x}) d\mathbf{x}.$$

Global probability of misclassification (global misclassification error):

$$e(r) = \sum_{k=1}^K p_k e_k(r).$$

Misclassification cost

Cost of misclassifying an observation of G_ℓ in G_k :

$$C : (k, \ell) \in \{1, \dots, K\} \times \{1, \dots, K\} \rightarrow C(k, \ell) \in \mathbb{R}^+,$$

with $C(k, k) = 0$.

Remarks:

- ▶ $C(k, \ell)$ generally not symmetric
- ▶ to be defined with practitioner (or fixed to 1 if you have no information)

Examples of misclassification cost

Medecine:

- ▶ a test allows to detect if a patient is ill (G_1) or not (G_2)
- ▶ $C(1, 2)$: cost of classifying a *healthy* patient (G_2) as *ill* (G_1) \Rightarrow cost of carrying out more exam on a healthy patient
- ▶ $C(2, 1)$: cost of classifying a *ill* patient as a *healthy* one \Rightarrow cost that an ill patient go back home without treatment

Finance:

- ▶ a consumer finance company predicts if the customer will have (G_1) problem in repaying loan or not (G_2)
- ▶ $C(1, 2)$: cost of classifying a *good* customer (G_2) as a *bad* one (G_1)
- ▶ $C(2, 1)$: cost of classifying a *bad* customer as a *good* one
- ▶ the company probably has to choose $C(2, 1) \gg C(1, 2)$

Bayes optimal classification rule

Conditional risk associated to \mathbf{x} : average cost of misclassification of \mathbf{x}

$$R(r(\mathbf{X})|\mathbf{X} = \mathbf{x}) = E[C(r(\mathbf{X}), Z)|\mathbf{X} = \mathbf{x}] = \sum_{k=1}^K C(r(\mathbf{x}), k)t_k(\mathbf{x}),$$

Average risk

$$R(r) = E_{\mathbf{X}}[R(r(\mathbf{X})|\mathbf{X} = \mathbf{x})] = \sum_{k=1}^K p_k \sum_{\ell=1}^K C(\ell, k) \int_{\Omega_{\ell}} f_k(\mathbf{x}) d\mathbf{x}.$$

Proofs: exercice.

Bayes optimal classification rule

We look for the optimal rule r^* which minimize the average risk, which is equivalent to minimize the conditional risk since:

$$R(r^*) = \min_r E_{\mathbf{X}}[R(r(\mathbf{X})|\mathbf{X} = \mathbf{x})] \geq E_{\mathbf{X}}[\min_r R(r(\mathbf{X})|\mathbf{X} = \mathbf{x})].$$

The optimal rule classifies \mathbf{x} into G_k if

$$R(r(\mathbf{X}) = k|\mathbf{X} = \mathbf{x}) < R(r(\mathbf{X}) = \ell|\mathbf{X} = \mathbf{x}) \quad \forall \ell \neq k.$$

Since

$$\begin{aligned} R(r(\mathbf{X}) = k|\mathbf{X} = \mathbf{x}) &= E[C(k, Z)|\mathbf{X} = \mathbf{x}] \\ &= \sum_{\ell=1}^K C(k, \ell)t_{\ell}(\mathbf{x}) = \sum_{\ell \neq k}^K C(k, \ell)t_{\ell}(\mathbf{x}), \end{aligned}$$

the **optimal Bayes classification rule** is:

$$r^*(\mathbf{x}) = k \quad \text{if} \quad \sum_{\ell \neq k}^K C(k, \ell)t_{\ell}(\mathbf{x}) < \sum_{\ell \neq k'}^K C(k', \ell)t_{\ell}(\mathbf{x}) \quad \forall k' \neq k.$$

Bayes optimal rule for equal costs

If $C(k, \ell) = c \forall k \neq \ell$, the conditional risk is

$$R(r(\mathbf{X}) = k | \mathbf{X} = \mathbf{x}) = c \sum_{\ell \neq k}^K t_{\ell}(\mathbf{x}) = c(1 - t_k(\mathbf{x})),$$

and thus $r^*(\mathbf{x}) = k$ if $c(1 - t_k(\mathbf{x})) < c(1 - t_{k'}(\mathbf{x})) \quad \forall k' \neq k$ or equivalently

$$r^*(\mathbf{x}) = k \quad \text{if } t_k(\mathbf{x}) > t_{k'}(\mathbf{x}) \quad \forall k' \neq k.$$

$\Rightarrow \mathbf{x}$ is classified into the group which has the greater posterior probability **maximum a posteriori**.

Bayes optimal rule for equal costs

If moreover $c = 1$, the average risk is

$$\begin{aligned}R(r) &= \sum_{k=1}^K p_k \sum_{\ell \neq k} \int_{\Omega_\ell} f_k(\mathbf{x}) d\mathbf{x} \\&= \sum_{k=1}^K p_k \int_{\bar{\Omega}_\ell} f_k(\mathbf{x}) d\mathbf{x} \\&= \sum_{k=1}^K p_k e_k(r) \\&= e(r)\end{aligned}$$

Bayes optimal rule for 2 groups

For 2 groups, we have

$$\begin{aligned} r^*(\mathbf{x}) &= 1 && \text{if } C(1, 2)t_2(\mathbf{x}) < C(2, 1)t_1(\mathbf{x}), \\ \text{and } r^*(\mathbf{x}) &= 2 && \text{if } C(2, 1)t_1(\mathbf{x}) < C(1, 2)t_2(\mathbf{x}), \end{aligned}$$

and by noting $g(\mathbf{x}) = \frac{C(2,1)t_1(\mathbf{x})}{C(1,2)t_2(\mathbf{x})}$, the Bayes optimal rule is:

$$\begin{aligned} r^*(\mathbf{x}) &= 1 && \text{if } g(\mathbf{x}) > 1, \\ \text{and } r^*(\mathbf{x}) &= 2 && \text{if } g(\mathbf{x}) < 1. \end{aligned}$$

$g(\mathbf{x}) = 1$ is the equation of the **separating surface**.

Continuous features: the Gaussian mixture

The Gaussian Mixture Model

The density of group k is

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_k)^t \Sigma_k^{-1} (\mathbf{x} - \mu_k)\right\}$$

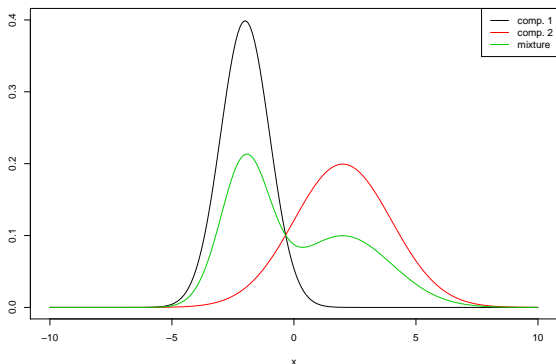
where

- ▶ μ_k is the mean vector
- ▶ Σ_k the covariance matrix of group k
- ▶ $|\Sigma_k|$ denotes the determinant of Σ_k

The Gaussian Mixture Model

An example of Gaussian mixture

```
x=seq(-10,10,.1)
plot(x,dnorm(x,-2,1),type='l',col=1,ylab='',ylim=c(0,0.4))
lines(x,dnorm(x,2,2),col=2)
lines(x,.5*dnorm(x,-2,1)+.5*dnorm(x,2,2),col=3)
legend('topright',legend=c('comp. 1','comp. 2','mixture'),
      col=1:3,lty=1)
```



Mixture model estimation in classification

Maximum likelihood estimation

- ▶ Estimation of the classification rule r^* is obtained by estimating θ by **maximum likelihood**
- ▶ likelihood in the classification context ($\underline{\mathbf{x}}, \underline{\mathbf{z}}$ available):

$$\begin{aligned} p(\underline{\mathbf{x}}, \underline{\mathbf{z}}) &= \prod_i p(\mathbf{x}_i, z_i) \\ &= \prod_i p(Z = z_i) f(\mathbf{x}_i | Z = z_i) \\ &= \prod_i \prod_k (p(Z = k) f(\mathbf{x}_i | Z = k))^{\tilde{z}_{ik}} \\ &= \prod_i \prod_k p_k^{\tilde{z}_{ik}} f_k(\mathbf{x}_i)^{\tilde{z}_{ik}} \end{aligned}$$

- ▶ Log-likelihood

$$\ell(\underline{\mathbf{x}}, \underline{\mathbf{z}}; \theta) = \sum_{i=1}^n \sum_{k=1}^K \tilde{z}_{ik} \left(\ln p_k - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_i - \mu_k)^t \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right)$$

Parameter estimation

Maximization leads to the usual empirical estimates:

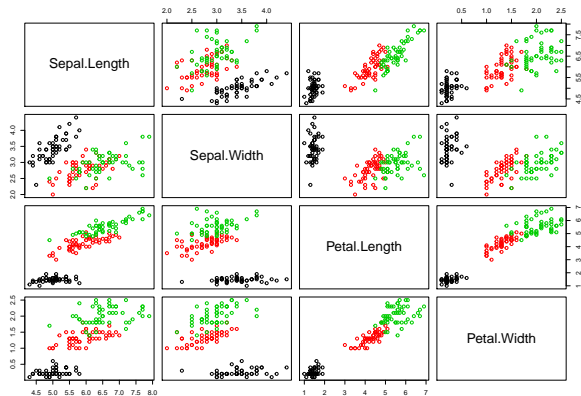
- ▶ $\hat{p}_k = \frac{n_k}{n}$ with $n_k = \sum_{i=1}^n \tilde{z}_{ik}$ the number of observations of group k
- ▶ $\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^n \tilde{z}_{ik} \mathbf{x}_i$
- ▶ $\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^n \tilde{z}_{ik} (\mathbf{x}_i - \hat{\mu}_k)^t (\mathbf{x}_i - \hat{\mu}_k)$

Exercise 1

Prove the expression of the previous estimators in the case of $p = 1$.
Use Lagrange multipliers for the constraint $\sum_{k=1}^K p_k = 1$.

Application on iris dataset

```
plot(iris[,1:4], col=iris$Species)
```



Application on iris dataset

Split into train and test data sets

```
set.seed(2)
tmp=sample(1:150,50)
X.train <- iris[-tmp,-5]
Class.train <- iris[-tmp,5]
X.test <- iris[tmp,-5]
Class.test <- iris[tmp,5]
```

Application on iris dataset

Estimation of a mixture of Gaussians (*model 'XXX' is the usual Gaussian model in MclustDA*)

```
irisMclustDA <- MclustDA(X.train, Class.train,  
                        modelName = "XXX")  
summary(irisMclustDA)
```

```
## -----  
## Gaussian finite mixture model for classification  
## -----  
##  
## MclustDA model summary:  
##  
##   log-likelihood    n df      BIC  
##      -113.9301 100 42 -421.2774  
##  
## Classes      n  % Model G  
##   setosa      35 35   XXX 1  
##   versicolor  35 35   XXX 1  
##   virginica   30 30   XXX 1  
##
```

Application on iris dataset

Evaluation of the prediction

```
tmp=summary(irisMclustDA, newdata = X.test,  
            newclass = Class.test)  
print(tmp$tab.newdata)
```

##		Predicted		
##	Class	setosa	versicolor	virginica
##	setosa	15	0	0
##	versicolor	0	14	1
##	virginica	0	0	20

Mixture model in classification

In comparison with other classification methods:

- ▶ MM has the advantage of interpretability
- ▶ but the classification power suffer from its assumption (each class should be Gaussian)

More flexibility can be introduce by considering mixture of mixture:

- ▶ each class can be itself a mixture

Mixture of mixture on iris dataset

Estimation of a mixture of Gaussians with selection by BIC of the number of mixture components per class (*option modelType = "MclustDA"*)

```
irisMclustDA <- MclustDA(X.train, Class.train,  
                        modelType = "MclustDA")  
summary(irisMclustDA)
```

```
## -----  
## Gaussian finite mixture model for classification  
## -----  
##  
## MclustDA model summary:  
##  
##   log-likelihood    n df      BIC  
##      -70.80528 100 63 -431.7363  
##  
## Classes      n  % Model G  
##   setosa      35 35   VEV 2  
##   versicolor 35 35   XXX 1  
##   virginica   30 30   VVE 2  
##
```

Application on iris dataset

Evaluation of the prediction (not necessary better for the simple iris data set)

```
tmp=summary(irisMclustDA, newdata = X.test,  
            newclass = Class.test)  
print(tmp$tab.newdata)
```

##		Predicted		
##	Class	setosa	versicolor	virginica
##	setosa	15	0	0
##	versicolor	0	14	1
##	virginica	0	0	20

Exercise 2

Implement your own maximum likelihood estimation for a Gaussian Mixture Model. Your function should be able to predict the class of a new observation.

Test it on simulated data.

Mixture model estimation in clustering

Maximum likelihood estimation

- ▶ Mixture model estimation for clustering is done by **maximum likelihood**
- ▶ likelihood in the clustering context (only $\underline{\mathbf{x}}$ available):

$$\begin{aligned} p(\underline{\mathbf{x}}) &= \prod_i p(\mathbf{x}_i) \\ &= \prod_i \sum_k p_k f_k(\mathbf{x}_i) \end{aligned}$$

- ▶ Log-likelihood

$$\ell(\underline{\mathbf{x}}; \boldsymbol{\theta}) = \sum_{i=1}^n \ln \left(\sum_{k=1}^K p_k \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^t \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)\right\} \right)$$

Not as easy to maximise as in the classification context !

The EM algorithm

The idea of Expectation-Maximization algorithm:

- ▶ to maximize the **completed**-likelihood (by the **unobserved data** \underline{z}) is easier than the observed-data likelihood:

$$\ell_c(\underline{\mathbf{x}}, \underline{\mathbf{z}}, \boldsymbol{\theta}) = \sum_{k=1}^K \sum_{i=1}^n \tilde{z}_{ik} \left(\ln p_k - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^t \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right)$$

- ▶ since \underline{z} is *unobserved*, the q -th iteration of the EM algorithm consists in alternating:
 - ▶ initialization: randomly choose $\boldsymbol{\theta}^{(0)}$
 - ▶ at iteration (q):
 - ▶ E step: computation of

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)}) = E[\ell_c(\underline{\mathbf{x}}, \underline{\mathbf{z}}, \boldsymbol{\theta}) | \underline{\mathbf{x}}, \boldsymbol{\theta}^{(q)}]$$

- ▶ M step: maximisation of $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})$ according to $\boldsymbol{\theta}$:

$$\boldsymbol{\theta}^{(q+1)} = \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})$$

until convergence of the log-likelihood: $|\ell(\underline{\mathbf{x}}, \boldsymbol{\theta}^{(q+1)}) - \ell(\underline{\mathbf{x}}, \boldsymbol{\theta}^{(q)})| < \epsilon$

EM algorithm - E step

Computation of $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)}) = E[\ell_c(\underline{\mathbf{x}}, \underline{\mathbf{z}}, \boldsymbol{\theta}) | \underline{\mathbf{x}}, \boldsymbol{\theta}^{(q)}]$:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)}) = \sum_{k=1}^K \sum_{i=1}^n E[\tilde{z}_{ik} | \underline{\mathbf{x}}, \boldsymbol{\theta}^{(q)}] \left(\ln p_k - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^t \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right)$$

with

$$\begin{aligned} E[\tilde{z}_{ik} | \underline{\mathbf{x}}, \boldsymbol{\theta}^{(q)}] &= 1 \times P(\tilde{z}_{ik} = 1 | \mathbf{x}_i, \boldsymbol{\theta}^{(q)}) + 0 \times P(\tilde{z}_{ik} = 0 | \mathbf{x}_i, \boldsymbol{\theta}^{(q)}) \\ &= \frac{f_{|\tilde{z}_{ik}=1}(\mathbf{x}_i, \boldsymbol{\theta}^{(q)}) P(\tilde{z}_{ik} = 1 | \boldsymbol{\theta}^{(q)})}{f(\mathbf{x}_i, \boldsymbol{\theta}^{(q)})} \\ &= \frac{f_k(\mathbf{x}_i, \boldsymbol{\theta}^{(q)}) p_k^{(q)}}{f_{\mathbf{X}}(\mathbf{x}_i, \boldsymbol{\theta}^{(q)})} \\ &= t_k^{(q)}(\mathbf{x}_i) \end{aligned}$$

EM algorithm - M step

Maximisation of $Q(\theta, \theta^{(q)})$ according to θ :

$$Q(\theta, \theta^{(q)}) = \sum_{k=1}^K \sum_{i=1}^n t_k^{(q)}(\mathbf{x}_i) \left(\ln p_k - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_i - \mu_k)^t \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right)$$

is equivalent to the log-likelihood maximization in the classification context, but by **ponderating** each observation by $t_k^{(q)}(\mathbf{x}_i)$

- ▶ $\hat{p}_k = \frac{n_k^{(q)}}{n}$ with $n_k = \sum_{i=1}^n t_k^{(q)}(\mathbf{x}_i)$
- ▶ $\hat{\mu}_k = \frac{1}{n_k^{(q)}} \sum_{i=1}^n t_k^{(q)}(\mathbf{x}_i) \mathbf{x}_i$
- ▶ $\hat{\Sigma}_k = \frac{1}{n_k^{(q)}} \sum_{i=1}^n t_k^{(q)}(\mathbf{x}_i) (\mathbf{x}_i - \mu_k)^t (\mathbf{x}_i - \mu_k)$

EM algorithm - properties

- ▶ the EM increases the likelihood at each step: \Rightarrow it converges to a local maximum of the likelihood
- ▶ convergence to the global maximum is expected to be achieved with multiple initializations of the algorithm
- ▶ in practice, the most efficient initialization strategy is:
 - ▶ run several small EM (with 10 iterations)
 - ▶ run a long EM starting from the small EM solution with highest log-likelihood

EM algorithm - proof of convergence

Since

$$p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})p(\mathbf{x}; \boldsymbol{\theta})$$

taking the logarithm we have:

$$\ell_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \ln p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}) + \ell(\mathbf{x}; \boldsymbol{\theta})$$

and then

$$\ell(\mathbf{x}; \boldsymbol{\theta}) = \ell_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) - \ln p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$$

Let's compute $E_{\theta^{(q)}}[\cdot|\mathbf{x}]$ of these terms:

$$\ell(\mathbf{x}; \boldsymbol{\theta}) = \underbrace{E_{\theta^{(q)}}[\ell_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})|\mathbf{x}]}_{Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})} - \underbrace{E_{\theta^{(q)}}[\ln p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})]}_{H(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})} \quad (1)$$

EM algorithm - proof of convergence

Let's look at $H(\theta, \theta^{(q)})$:

$$\begin{aligned}H(\theta^{(q)}, \theta^{(q)}) - H(\theta, \theta^{(q)}) &= E_{\theta^{(q)}}[\ln p(\mathbf{z}|\mathbf{x}; \theta^{(q)}) - \ln p(\mathbf{z}|\mathbf{x}; \theta)] \\&= E_{\theta^{(q)}}\left[\ln \frac{p(\mathbf{z}|\mathbf{x}; \theta^{(q)})}{p(\mathbf{z}|\mathbf{x}; \theta)}\right] \\&= \int \ln \frac{p(\mathbf{z}|\mathbf{x}; \theta^{(q)})}{p(\mathbf{z}|\mathbf{x}; \theta)} p(\mathbf{z}|\mathbf{x}; \theta^{(q)}) d\mathbf{z} \\&= KL(p(\mathbf{z}|\mathbf{x}; \theta^{(q)}), p(\mathbf{z}|\mathbf{x}; \theta)) \\&\geq 0\end{aligned}$$

Consequently, for all θ

$$H(\theta, \theta^{(q)}) \leq H(\theta^{(q)}, \theta^{(q)})$$

EM algorithm - proof of convergence

Since at each M step, we look for

$$\boldsymbol{\theta}^{(q+1)} = \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})$$

we have:

$$Q(\boldsymbol{\theta}^{(q+1)}, \boldsymbol{\theta}^{(q)}) \geq Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})$$

for all $\boldsymbol{\theta}$ and in particular for $\boldsymbol{\theta} = \boldsymbol{\theta}^{(q)}$.

Consequently:

$$Q(\boldsymbol{\theta}^{(q+1)}, \boldsymbol{\theta}^{(q)}) \geq Q(\boldsymbol{\theta}^{(q)}, \boldsymbol{\theta}^{(q)})$$

EM algorithm - proof of convergence

Using (1) with $\theta = \theta^{(q+1)}$ we have

$$\ell(\mathbf{x}; \theta^{(q+1)}) = \underbrace{Q(\theta^{(q+1)}, \theta^{(q)})}_{\geq Q(\theta^{(q)}, \theta^{(q)})} - \underbrace{H(\theta^{(q+1)}, \theta^{(q)})}_{\leq H(\theta^{(q)}, \theta^{(q)})}$$

and thus

$$\begin{aligned}\ell(\mathbf{x}; \theta^{(q+1)}) &\geq Q(\theta^{(q)}, \theta^{(q)}) - H(\theta^{(q)}, \theta^{(q)}) \\ &\geq \ell(\mathbf{x}; \theta^{(q)})\end{aligned}$$

\Rightarrow After each Mstep of the EM algorithm, the likelihood increases

Exercise 3

Implement your own EM algorithm for Gaussian mixture model estimation.

Test it on simulated data.

Compare multiple random initialization with kmeans initialization.

Some variants of the EM algorithm

The **Classification EM** algorithm is a variant of the EM algorithm, obtained by *rounding* the $t_k^{(q)}(\mathbf{x}_i)$:

- ▶ $t_k^{(q)}(\mathbf{x}_i) = 1$ for the group k s.t. $t_k^{(q)}(\mathbf{x}_i)$ is maximum
- ▶ $t_k^{(q)}(\mathbf{x}_i) = 0$ for the other groups

CEM properties:

- ▶ CEM performs **hard** classification whereas EM performs **soft** classification
- ▶ the convergence of CEM is faster than EM, but leads to a biased estimator
- ▶ nevertheless, for large samples and well separated groups, the CEM is very efficient
- ▶ CEM can be a good way to initialize an EM

Some variants of the EM algorithm

The **Stochastic EM** algorithm is a variant of the EM algorithm, obtained by *generating* the z_i according to the probabilities $t_k^{(q)}(\mathbf{x}_i)$.

- ▶ after a burn-in period, SEM generate a sample of $\theta^{(q)}$ whose distribution is *around* the maximum likelihood
- ▶ final estimation can be obtained by the mean/median of this generated distribution
- ▶ SEM could be an alternative to EM for more complex model in which the $t_k^{(q)}(\mathbf{x}_i)$ are intractable

A particular Gaussian mixture model

Mixture models sometimes generalized well known clustering algorithm

- ▶ assume (for parcimony) that $\Sigma_k = \alpha I_p$ for every clusters
- ▶ assume equal proportions: $\pi_1 = \dots = \pi_K$
- ▶ estimate the model with the CEM algorithm

A particular Gaussian mixture model

Mixture models sometimes generalized well known clustering algorithm

- ▶ assume (for parcimony) that $\Sigma_k = \alpha I_p$ for every clusters
- ▶ assume equal proportions: $\pi_1 = \dots = \pi_K$
- ▶ estimate the model with the CEM algorithm

⇒ you are running **the k-means algorithm**

Exercise 4

Implement your own kmeans algorithm.

Test it on simulated data.