

# Session 1: Bayesian Decision theory & Mixture Model

HMM for TS  
classif &  
filtering

Stéphane Derrode, École Centrale de Lyon  
stephane.derrode@ec-lyon.fr

# Outline of session 1 (8h)

2

## 3. Mixture Model

- Definition, simulations.
- Automatic parameters learning: EM principle
- Lab: Gaussian mixture model and Image processing

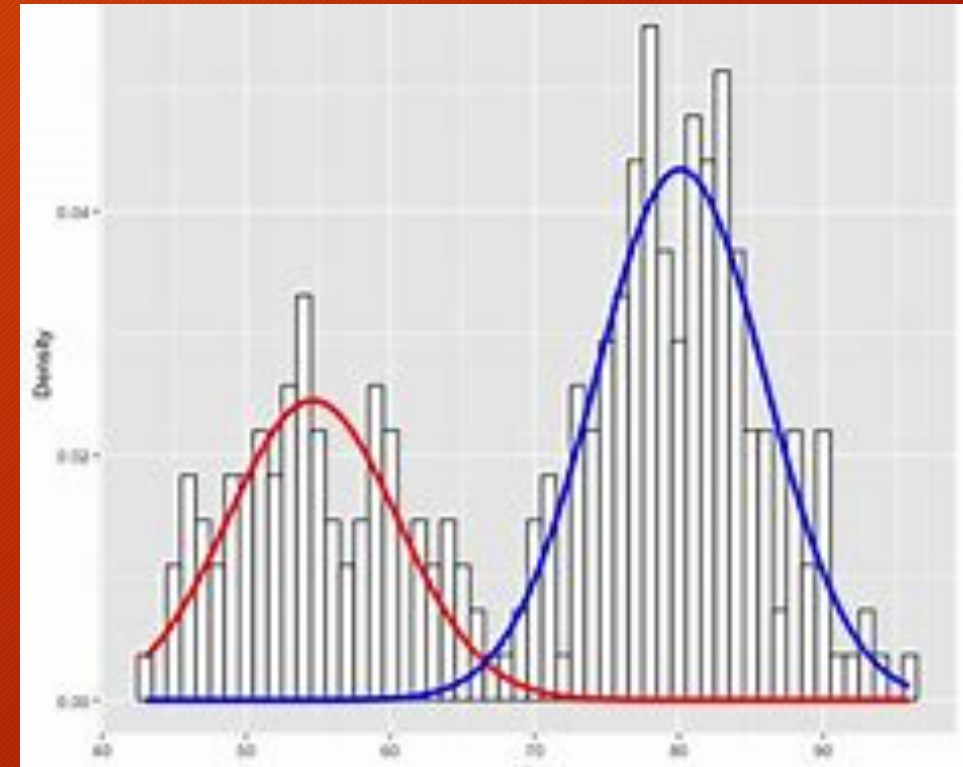
## 4. 2D Mixture Model

- 2D Gaussians
  - 2D mixture and EM re-estimation update
- 
- References and « Not seen! »
  - 1h sitting-exam, date: October 1, 2019

# 3. Mixture Model

3

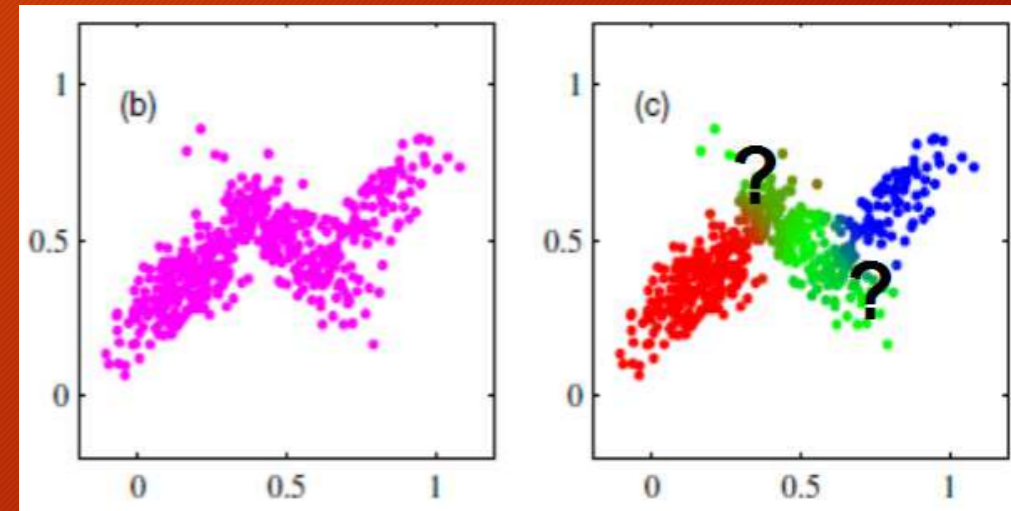
- Definition, simulations
- Exemple of mixture : image processing
- Automatic parameters learning: the EM principle
- Exercise : Gaussian mixture model



# 3. Mixture Model

In statistics, a mixture model is a probabilistic model for representing the presence of subpopulations within an overall population, without requiring that an observed data set should identify the sub-population to which an individual observation belongs.

Formally a mixture model corresponds to the mixture distribution that represents the probability distribution of observations in the overall population. However, while problems associated with "mixture distributions" relate to deriving the properties of the overall population from those of the sub-populations, "mixture models" are used to make statistical inferences about the properties of the sub-populations given only observations on the pooled population, without sub-population identity information.



# 3. Mixture Model

5

Suppose that we have a sample

$$\mathbf{y} = \mathbf{y}_1^N = \{y_1, y_2, \dots, y_n, \dots, y_N\}$$

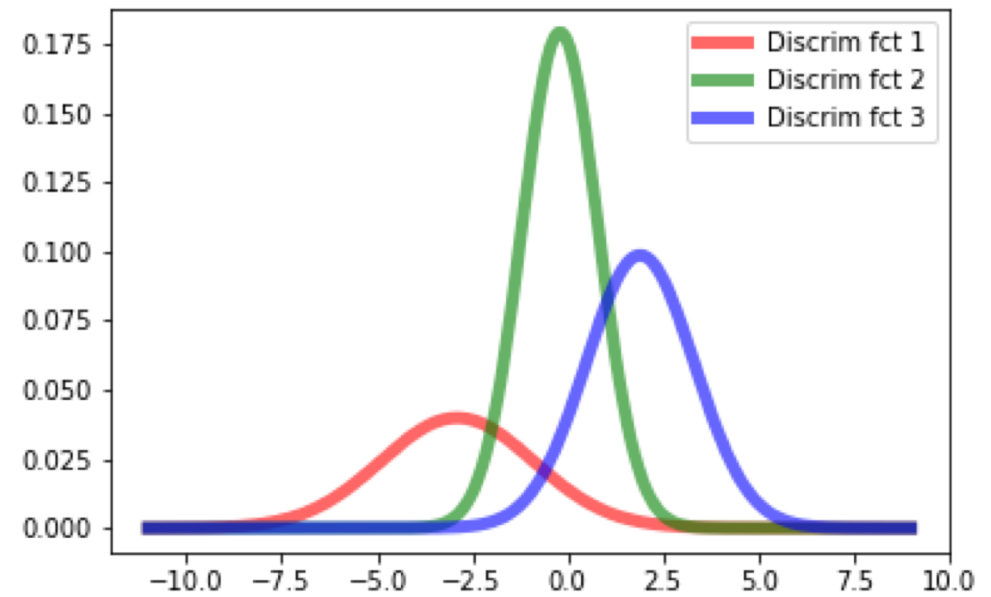
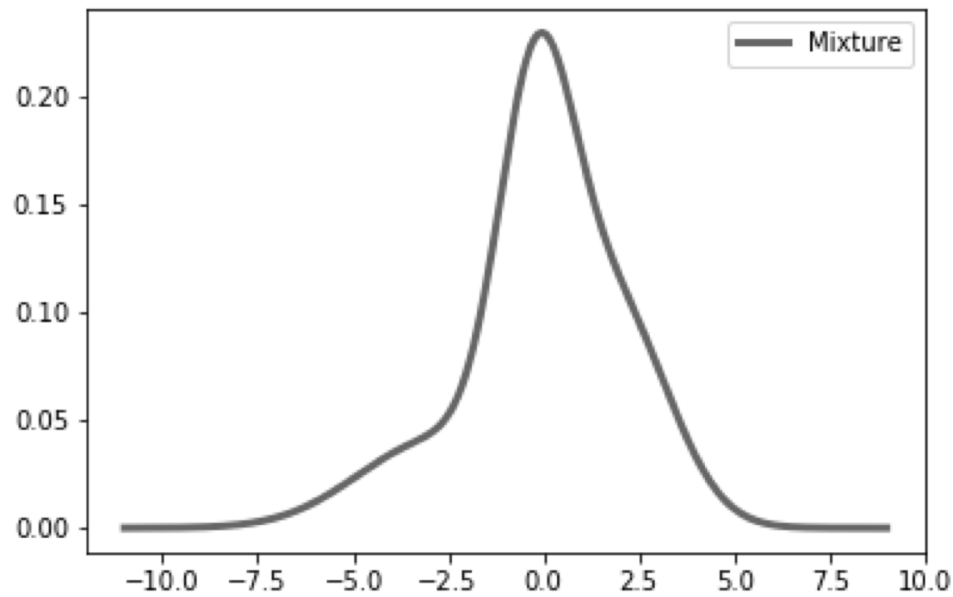
distributed according to a mixture of Gaussian distributions, so that all sample has the following with density:

$$P(Y_n = y_n) = f(y_n) = \sum_{k=1}^K \pi_k f_k(y_n)$$

A Gaussian mixture model is made with Gaussian  $f_k$ .

# 3. Mixture Model

6



$$\mathcal{N}(\mu_1 = -2.9, \sigma_1 = 2)$$

$$\pi_1 = 0.10$$

$$\mathcal{N}(\mu_2 = -0.2, \sigma_2 = 1)$$

$$\pi_2 = 0.55$$

$$\mathcal{N}(\mu_3 = 1.9, \sigma_3 = \sqrt{2})$$

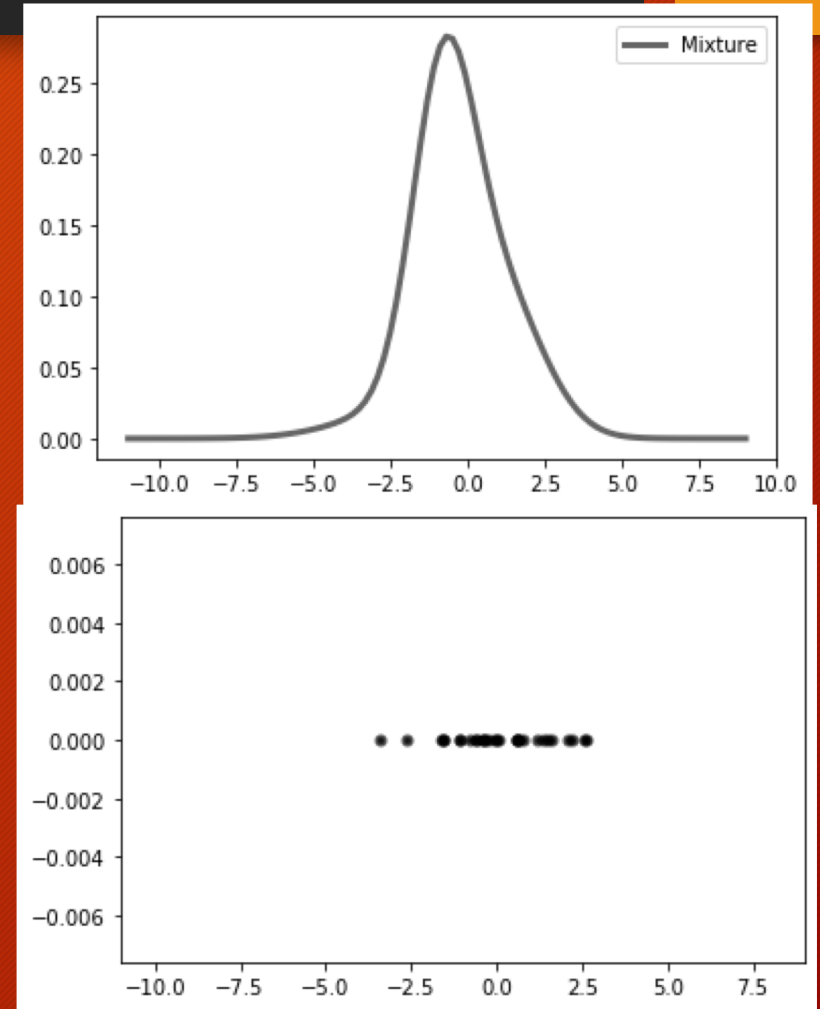
$$\pi_3 = 0.35$$

# 3. Mixture Model

Question : How to draw a sample for a mixture ?

$$P(Y_n = y_n) = f(y_n) = \sum_{k=1}^K \pi_k f_k(y_n)$$

$$\begin{aligned} \mathcal{N}(\mu_1 = -2.9, \sigma_1 = 2) & \quad \pi_1 = 0.10 \\ \mathcal{N}(\mu_2 = -0.2, \sigma_2 = 1) & \quad \pi_2 = 0.55 \\ \mathcal{N}(\mu_3 = 1.9, \sigma_3 = \sqrt{2}) & \quad \pi_3 = 0.35 \end{aligned}$$

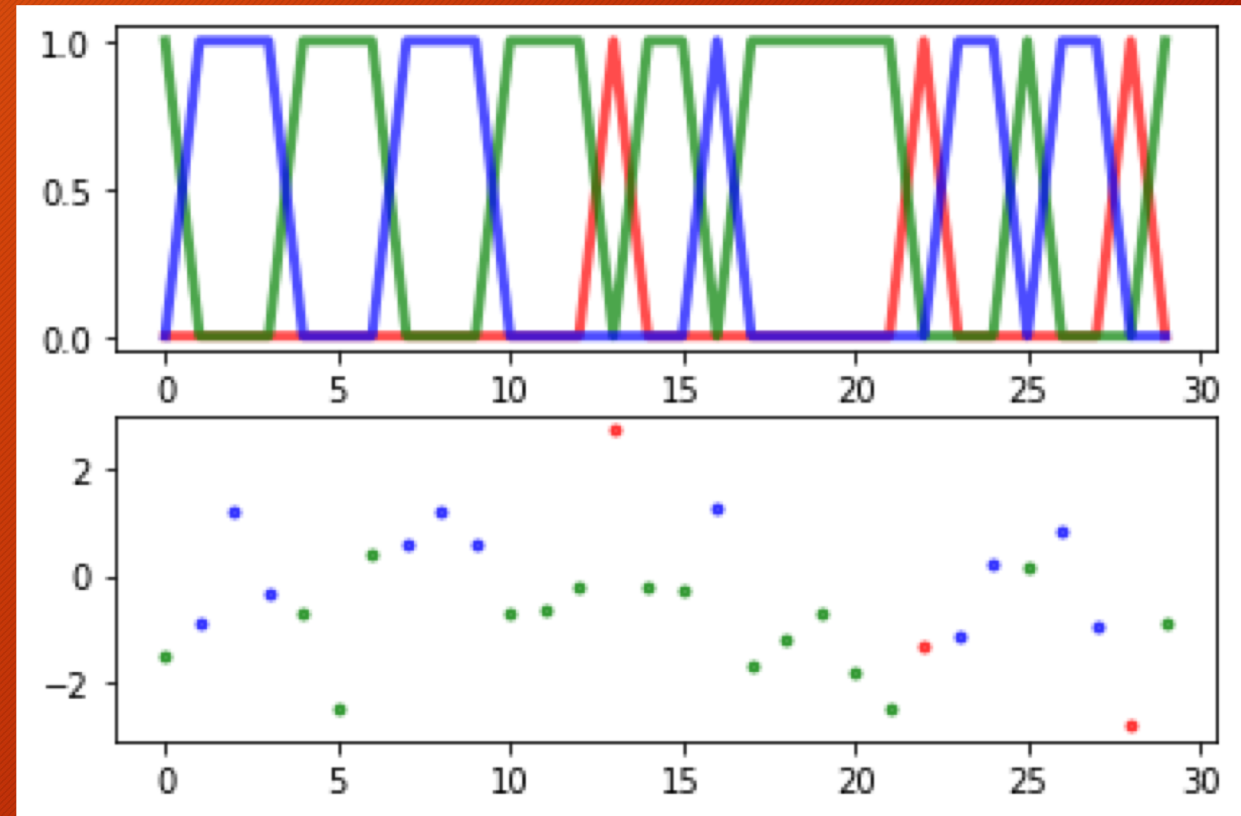


# 3. Mixture Model

Question : How to draw a sample for a mixture ?

1. Sampling according to the a priori proba to get the class number.
2. Sampling according to the selected Gaussian.

$$\begin{aligned} \mathcal{N}(\mu_1 = -2.9, \sigma_1 = 2) & \quad \pi_1 = 0.10 \\ \mathcal{N}(\mu_2 = -0.2, \sigma_2 = 1) & \quad \pi_2 = 0.55 \\ \mathcal{N}(\mu_3 = 1.9, \sigma_3 = \sqrt{2}) & \quad \pi_3 = 0.35 \end{aligned}$$



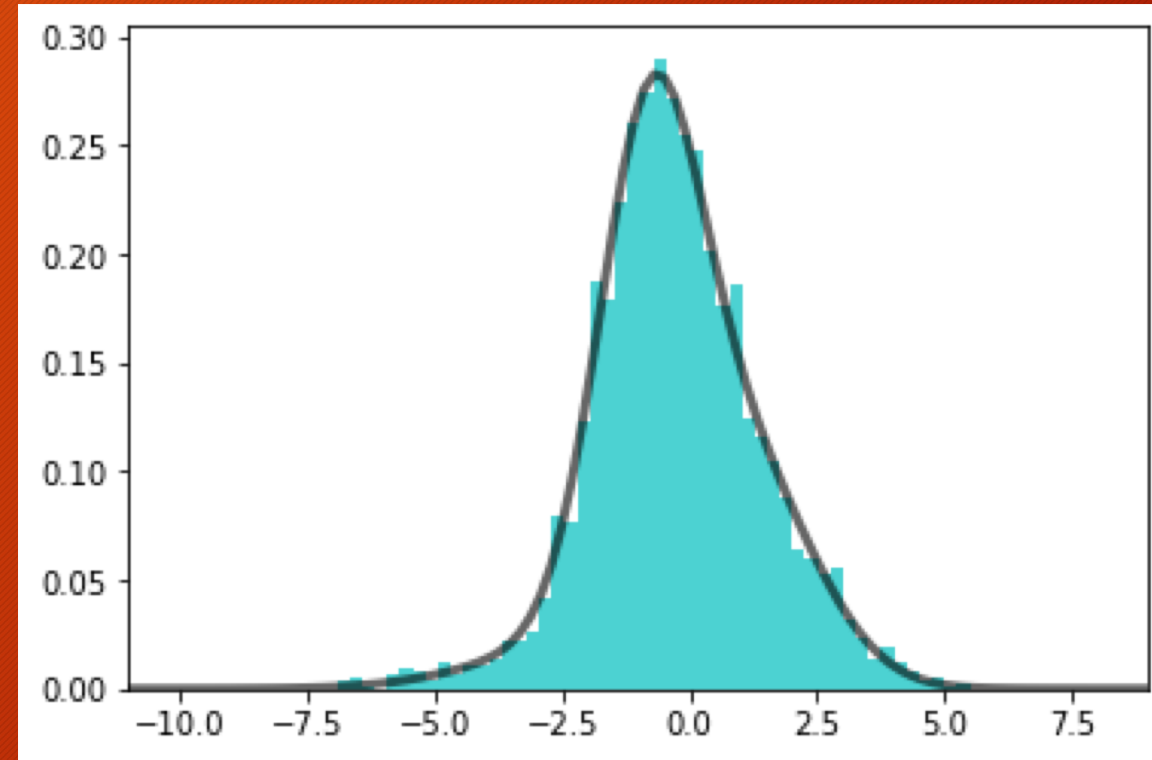


# 3. Mixture Model

Question : How to draw a sample for a mixture ?

1. Sampling according to the a priori proba to get the class number.
2. Sampling according to the selected Gaussian.

$$\begin{aligned} \mathcal{N}(\mu_1 = -2.9, \sigma_1 = 2) & \quad \pi_1 = 0.10 \\ \mathcal{N}(\mu_2 = -0.2, \sigma_2 = 1) & \quad \pi_2 = 0.55 \\ \mathcal{N}(\mu_3 = 1.9, \sigma_3 = \sqrt{2}) & \quad \pi_3 = 0.35 \end{aligned}$$



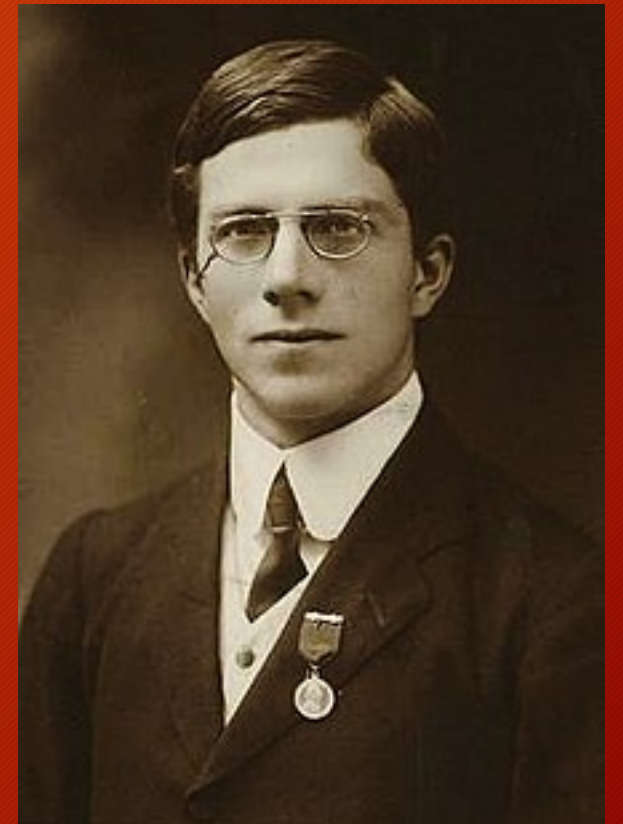
For a sample 3000 data

# 3. Mixture Model

In statistics, the likelihood expresses how probable a given set of observations is for different values of statistical parameters. It is equal to the joint probability distribution of the random sample evaluated at the given observations, and it is, thus, solely a function of parameters that index the family of those probability distributions.

For MM 
$$\mathcal{L}_{\Theta}(\mathbf{y}) = \prod_{n=1}^N \sum_{k=1}^K \pi_k f_k(y_n)$$

considered as a function of  $\Theta$ .



Sir Ronald Fisher

# 3. Mixture Model

11

Mapping from the parameter space to the real line, the likelihood function presents a peak, if it exists, which represents the combination of model parameter values that maximize the probability of drawing the sample actually obtained.

The procedure for obtaining these arguments of the maximum of the likelihood function is known as maximum likelihood estimation, which for computational convenience is usually done using the natural logarithm of the likelihood, known as the log-likelihood function.

Question: How to maximise  $\mathcal{L}_{\Theta}(\mathbf{y})$ ?

# 3. Mixture Model

Question: How to maximise  $\mathcal{L}_\Theta(\mathbf{y})$  ?

- Direct maximization is not possible.
- Solution : Expectation-Maximization (EM) algorithm

We define the `joint likelihood`

$$\mathcal{H}_\Theta(\mathbf{y}, \mathbf{X}) = \prod_{n=1}^N \pi_{X_n} f_{X_n}(y_n) = \prod_{n=1}^N \sum_{k=1}^K \pi_k f_k(y_n) \mathbb{I}_{(X_n=k)}$$

This a random function of  $\mathbf{X} = \mathbf{X}_1^N = \{X_1, X_2, \dots, X_n, \dots, X_N\}$

# 3. Mixture Model

The EM algorithm: this is an iterative algorithm to estimate the maximum of the likelihood function, by computing iteratively two steps:

1. Expectation of the auxiliary function

where 
$$Q\left(\Theta; \Theta^{(\ell)}\right) = E\left[\ln \mathcal{H}_{\Theta}(\mathbf{y}, \mathbf{X}) \mid \Theta^{(\ell)}\right]$$

- $\Theta$  is the set of true parameters (we are looking for).
- $\Theta^{(\ell)}$  is the estimated parameters set at iteration  $\ell$ .

2. Maximization of the auxiliary function

$$\Theta^{(\ell+1)} = \arg \max_{\Theta} Q\left(\Theta; \Theta^{(\ell)}\right)$$

## 3. Mixture Model

14

Properties of the EM algorithm (not proven)

1. Construction of a series of estimators for which the likelihood is increasing.

$$\mathcal{L}_{\Theta^{(e+1)}}(\mathbf{y}) \geq \mathcal{L}_{\Theta^{(e)}}(\mathbf{y})$$

The likelihood is always increasing (this is a sufficient condition to ensure the convergence of the EM algorithm).

# 3. Mixture Model

Properties of the EM algorithm (not proven)

- 2. Convergence towards one of the (local) maxima of likelihood since we have

$$\left. \frac{\partial Q(\Theta; \Theta^{(\ell)})}{\partial \Theta} \right|_{\Theta = \Theta^{(\ell)}} = \left. \frac{\partial \mathcal{L}_{\Theta}(\mathbf{y})}{\partial \Theta} \right|_{\Theta = \Theta^{(\ell)}}$$

Initialization : Biernacki, C., Celeux, G. and Govaert, G. (2003). *Choosing starting values for the EM algorithm for getting the highest likelihood in multivariate Gaussian mixture models*. Computational Statistics and data analysis 41, 561-575.

# 3. Mixture Model

The joint log-likelihood is written

$$\ln \mathcal{H}_{\Theta}(\mathbf{y}, \mathbf{X}) = \sum_{n=1}^N \sum_{k=1}^K \ln (\pi_k f_k(y_n)) \mathbb{I}_{(X_n=k)}$$

At iteration  $\ell$ , the Gaussian pdfs write  $f_k^{(\ell)}$  with parameters  $\mu_k^{(\ell)}, \sigma_k^{(\ell)}$  and the auxiliary function writes

$$Q(\Theta; \Theta^{(\ell)}) = \sum_{n=1}^N \sum_{k=1}^K \ln \left( \pi_k^{(\ell)} f_k^{(\ell)}(y_n) \right) E \left[ \mathbb{I}_{(X_n=k)} | \mathbf{y}, \Theta^{(\ell)} \right]$$



# 3. Mixture Model

$$Q\left(\Theta; \Theta^{(\ell)}\right) = \sum_{n=1}^N \sum_{k=1}^K \ln\left(\pi_k^{(\ell)} f_k^{(\ell)}(y_n)\right) E\left[\mathbb{I}_{(X_n=k)} \mid \mathbf{y}, \Theta^{(\ell)}\right]$$

with

$$\begin{aligned} E\left[\mathbb{I}_{(X_n=k)} \mid \mathbf{y}, \Theta^{(\ell)}\right] &= p\left(X_n = k \mid \mathbf{y}, \Theta^{(\ell)}\right) = p\left(X_n = k \mid y_n, \Theta^{(\ell)}\right) \\ &= \gamma_n^{(\ell)}(k) = \frac{\pi_k^{(\ell)} f_k^{(\ell)}(y_n)}{\sum_{j=1}^K \pi_j^{(\ell)} f_j^{(\ell)}(y_n)} \end{aligned}$$

# 3. Mixture Model

Gaussian mixture: as an exercise, proof that EM-based re-estimation formulas for parameters of a MM can be written:

$$\pi_k^{(\ell+1)} = \frac{1}{N} \sum_{n=1}^N \gamma_n^{(\ell)}(k)$$

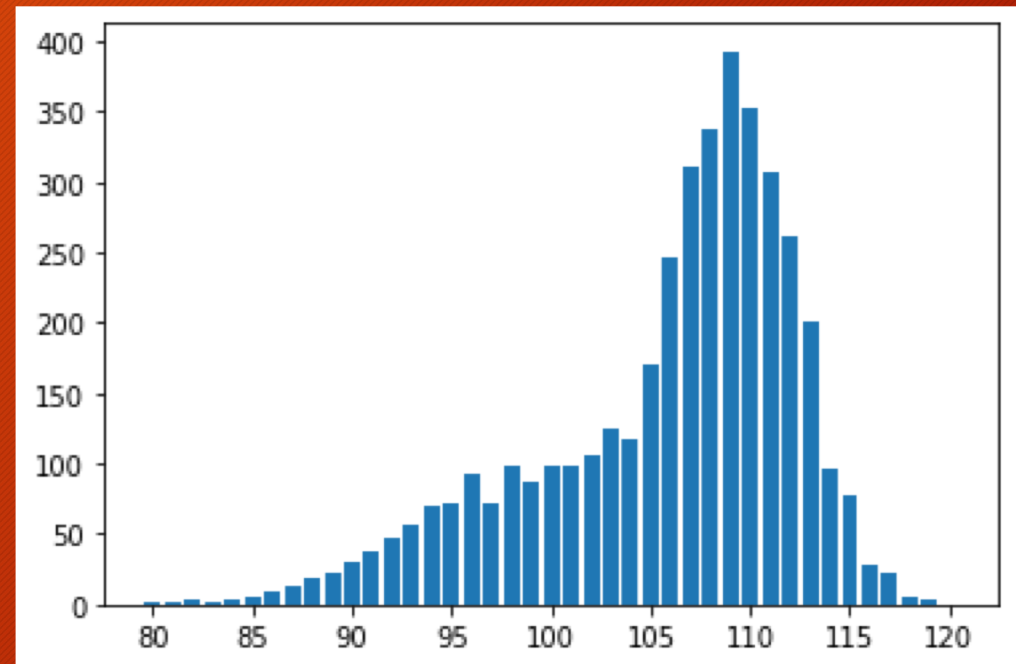
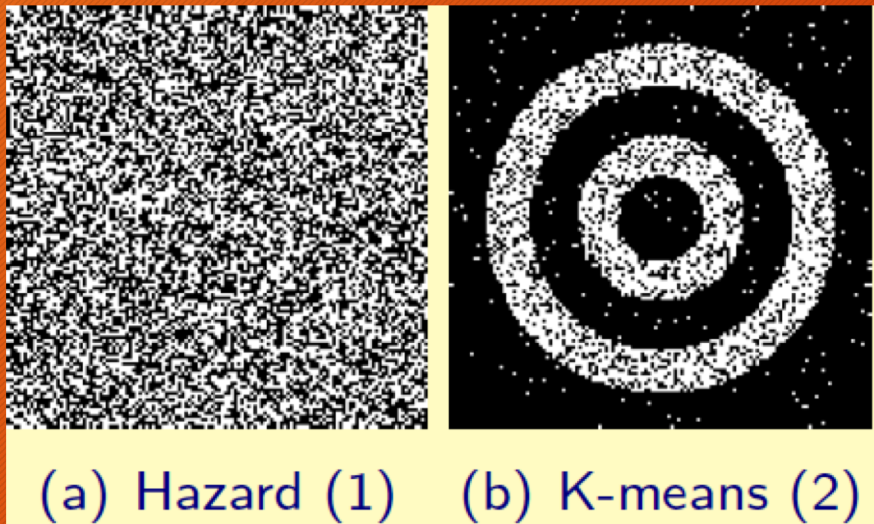
$$\mu_k^{(\ell+1)} = \frac{\sum_{n=1}^N \gamma_n^{(\ell)}(k) y_n}{\sum_{n=1}^N \gamma_n^{(\ell)}(k)}$$

$$\sigma_k^{2,(\ell+1)} = \frac{\sum_{n=1}^N \gamma_n^{(\ell)}(k) \left( y_n - \mu_k^{(\ell+1)} \right)^2}{\sum_{n=1}^N \gamma_n^{(\ell)}(k)}$$

# 3. Mixture Model

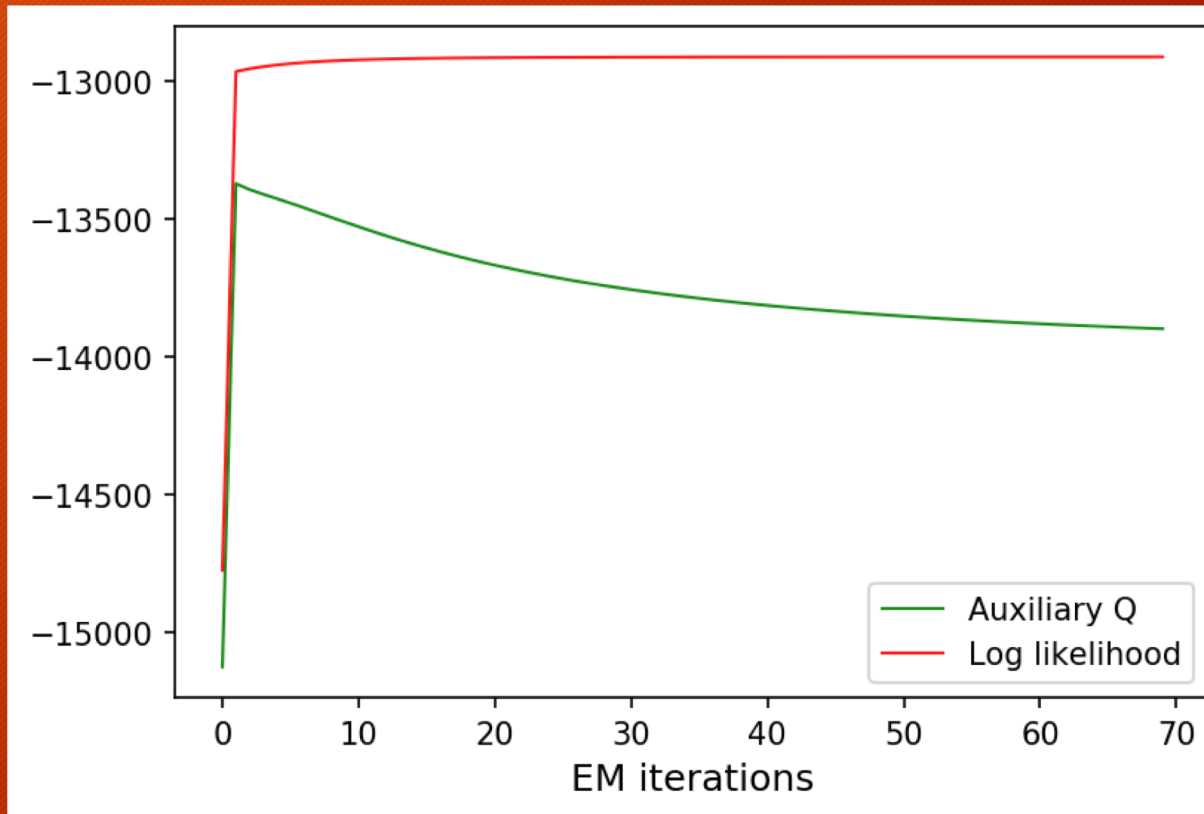
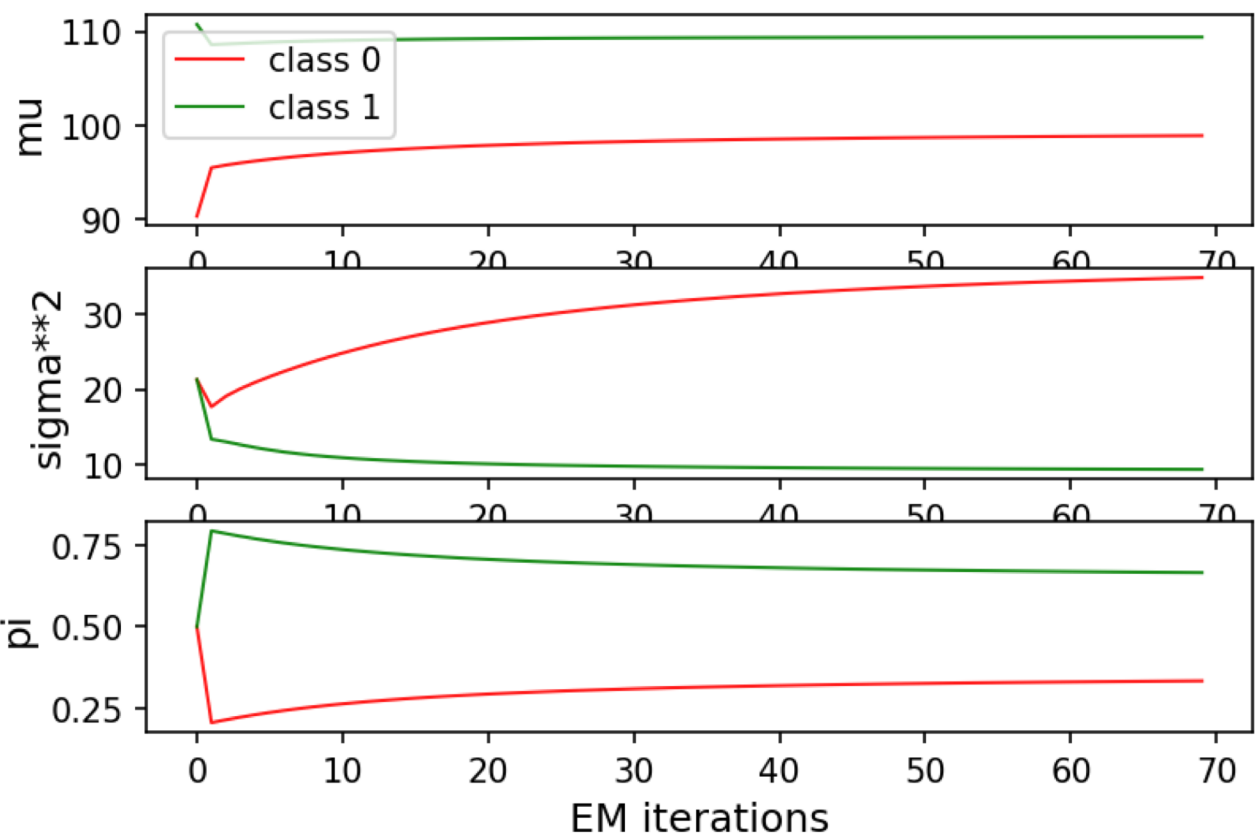
The EM algorithm looks for a local maxima of the likelihood: it requires the parameters to be initialized “not so far” from the true values.

Any idea to initialize parameters ?



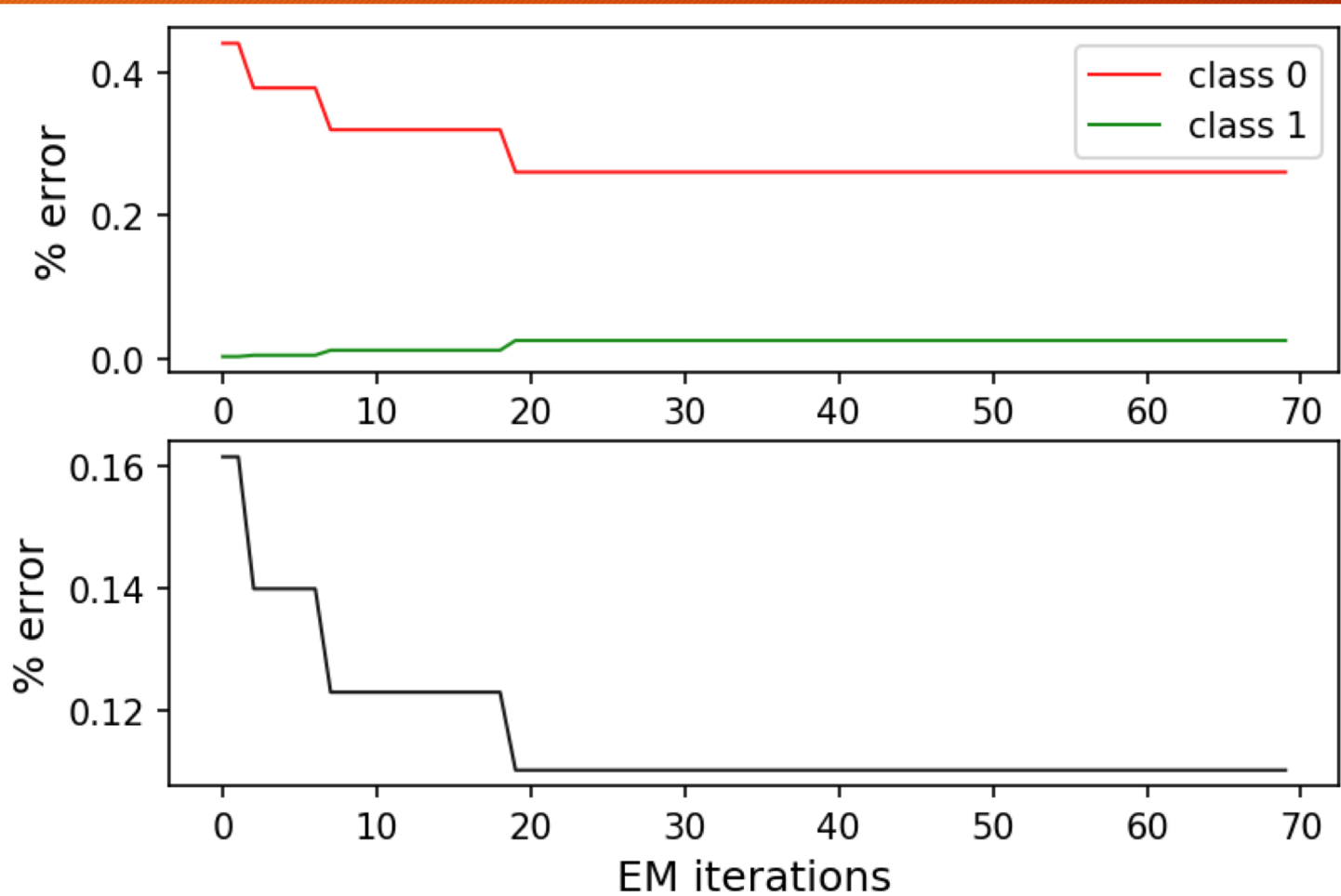
# 3. Mixture Model

20



# 3. Mixture Model

21



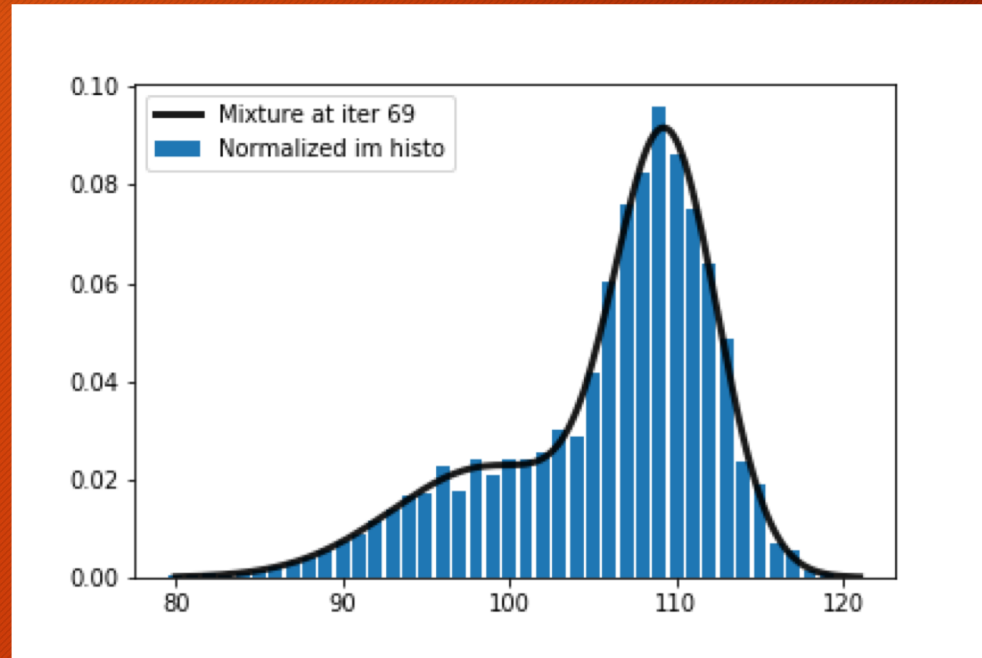
# 3. Mixture Model



$$\xi_1 = 0.261$$

$$\xi_2 = 0.025$$

$$\xi = 0.11$$



$$\pi_1 = 0.334 \quad \mathcal{N}(\mu_1 = 98.85, \sigma_1 = 5.91)$$

$$\pi_2 = 0.666 \quad \mathcal{N}(\mu_2 = 109.39, \sigma_2 = 3.06)$$

## 4. 2D Mixture Model

23

$$f(\mathbf{y}) = \sum_{k=1}^K \pi_k f_k(\mathbf{y}) \longrightarrow \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

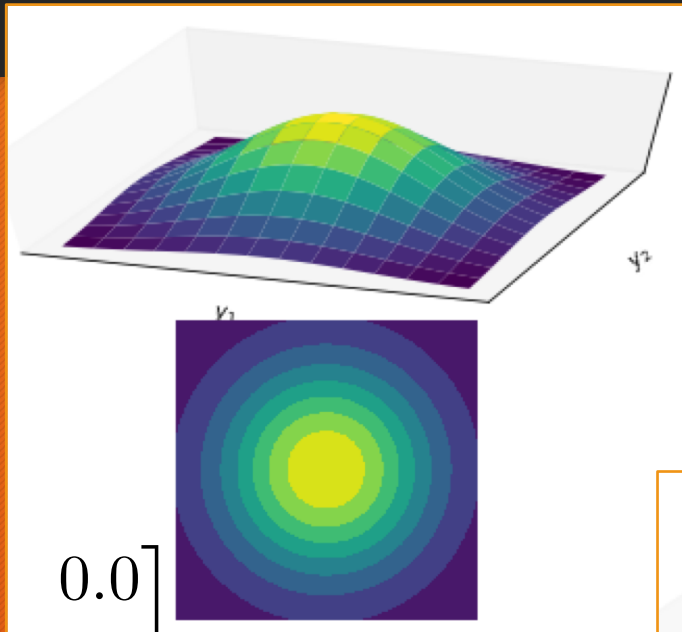
2D Gaussian

$$f(\mathbf{y}) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} e^{-\frac{1}{2}(\mathbf{y}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}-\boldsymbol{\mu})}$$

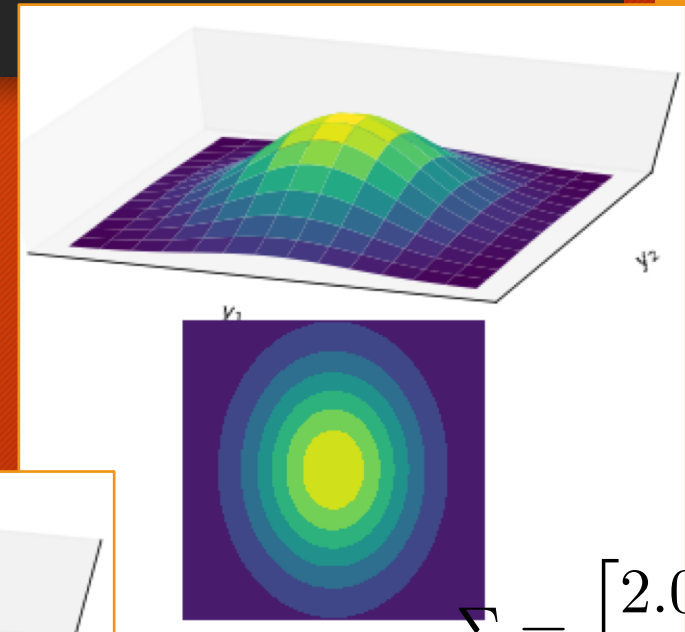
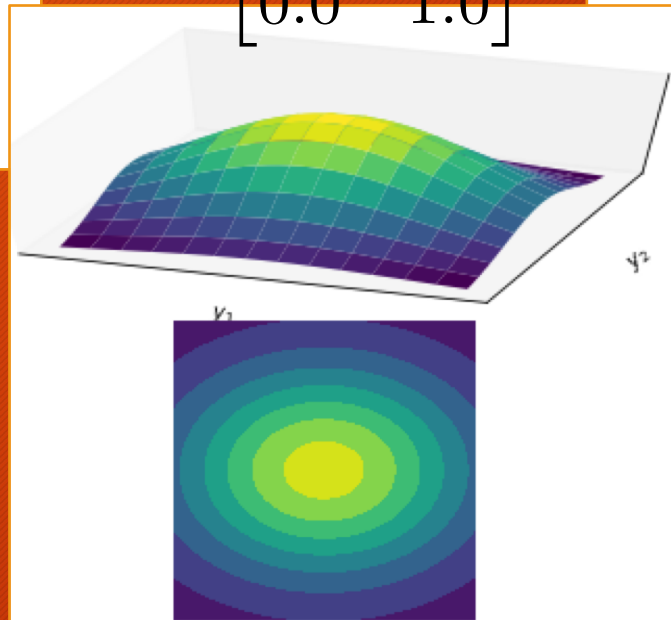
$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

# 4. 2D Mixture Model

$$\Sigma = \begin{bmatrix} 1.0 & 0.0 \\ 0.0 & 1.0 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 0.6 & 0.0 \\ 0.0 & 1.0 \end{bmatrix}$$

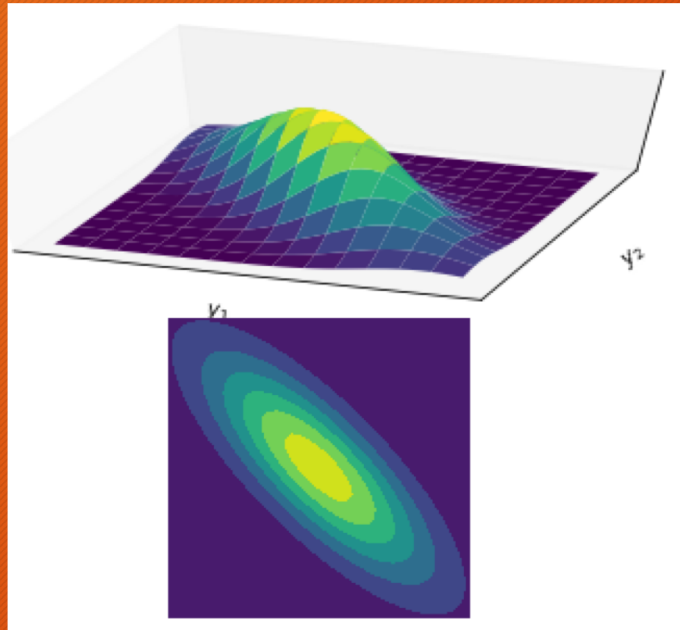


$$\Sigma = \begin{bmatrix} 2.0 & 0.0 \\ 0.0 & 1.0 \end{bmatrix}$$

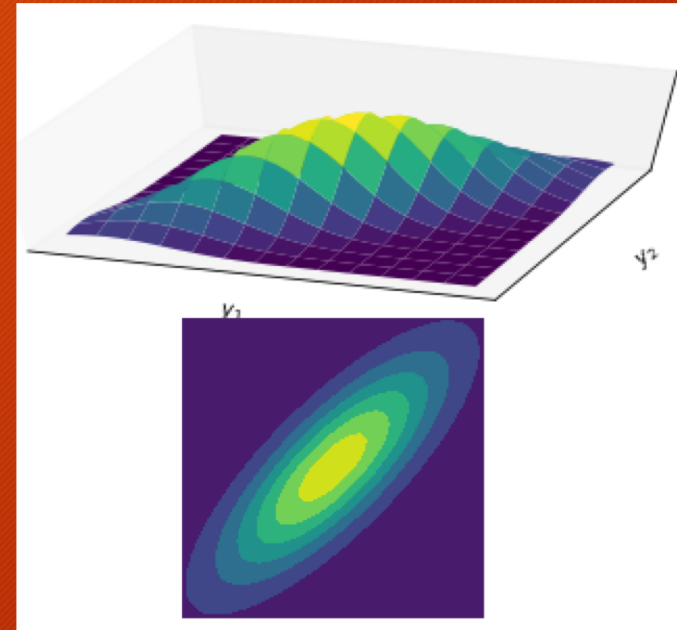


# 4. 2D Mixture Model

25



$$\Sigma = \begin{bmatrix} 1.0 & -0.8 \\ -0.8 & 1.0 \end{bmatrix}$$

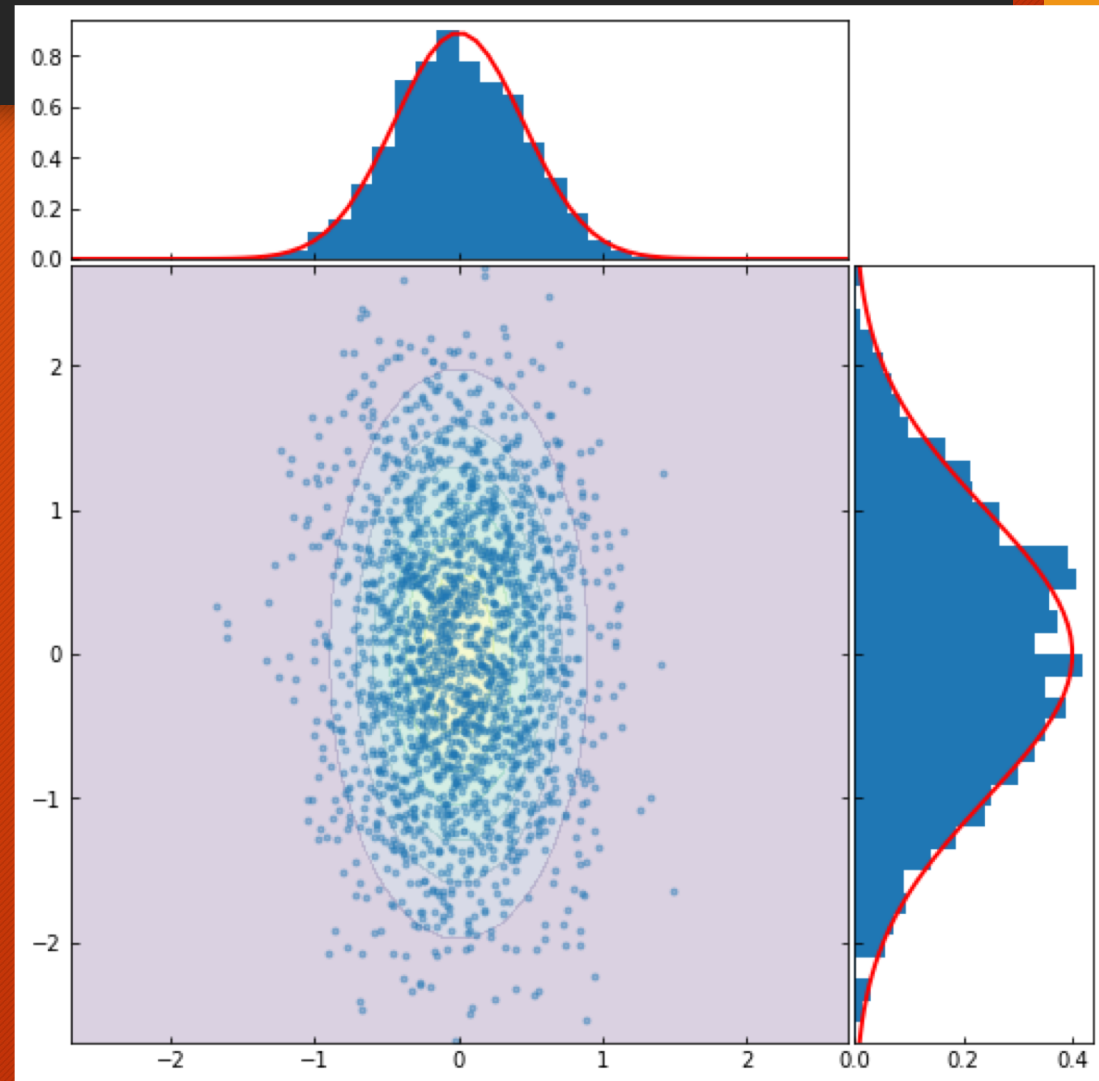


$$\Sigma = \begin{bmatrix} 1.0 & 0.8 \\ 0.8 & 1.0 \end{bmatrix}$$

## 4. 2D Mixture Model

26

Margins, conditional laws,  
empirical estimation of  
parameters



# 4. 2D Mixture Model

27

$$f(\mathbf{y}) = \sum_{k=1}^2 \pi_k f_k(\mathbf{y})$$

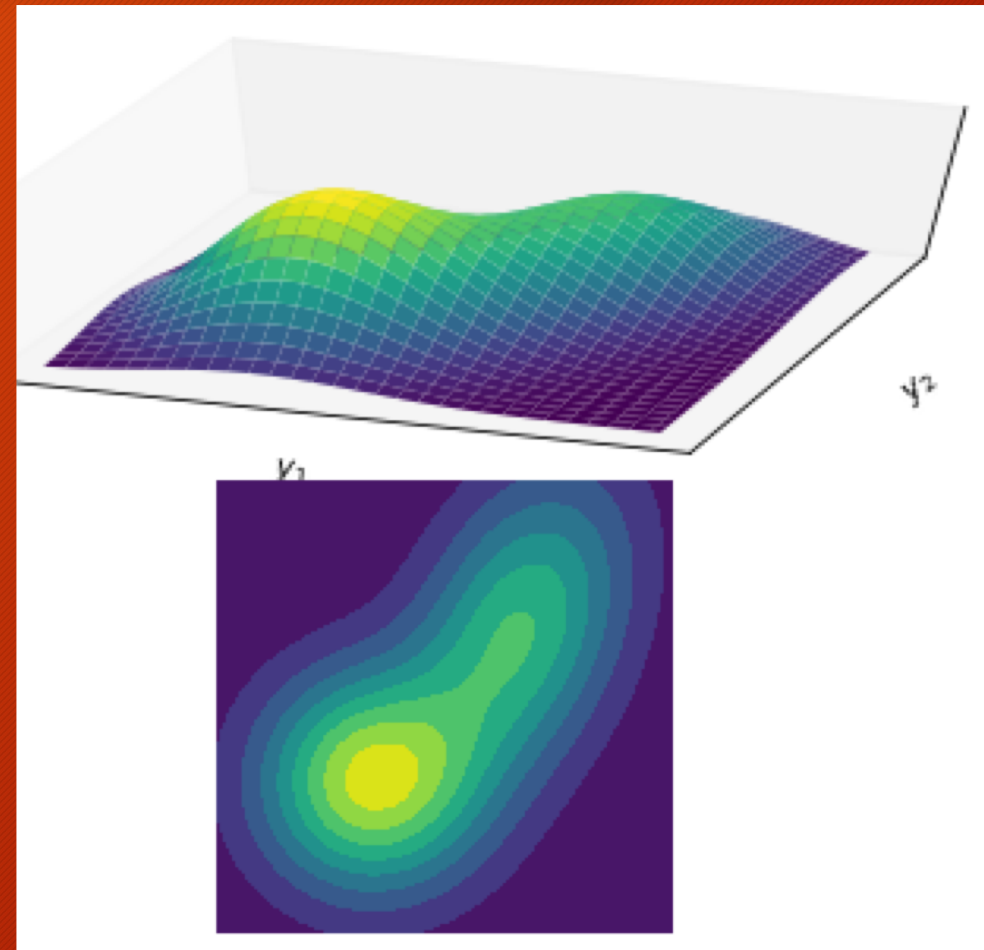
$$\pi_1 = \pi_2 = \frac{1}{2}$$

$$\mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\mu_2 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

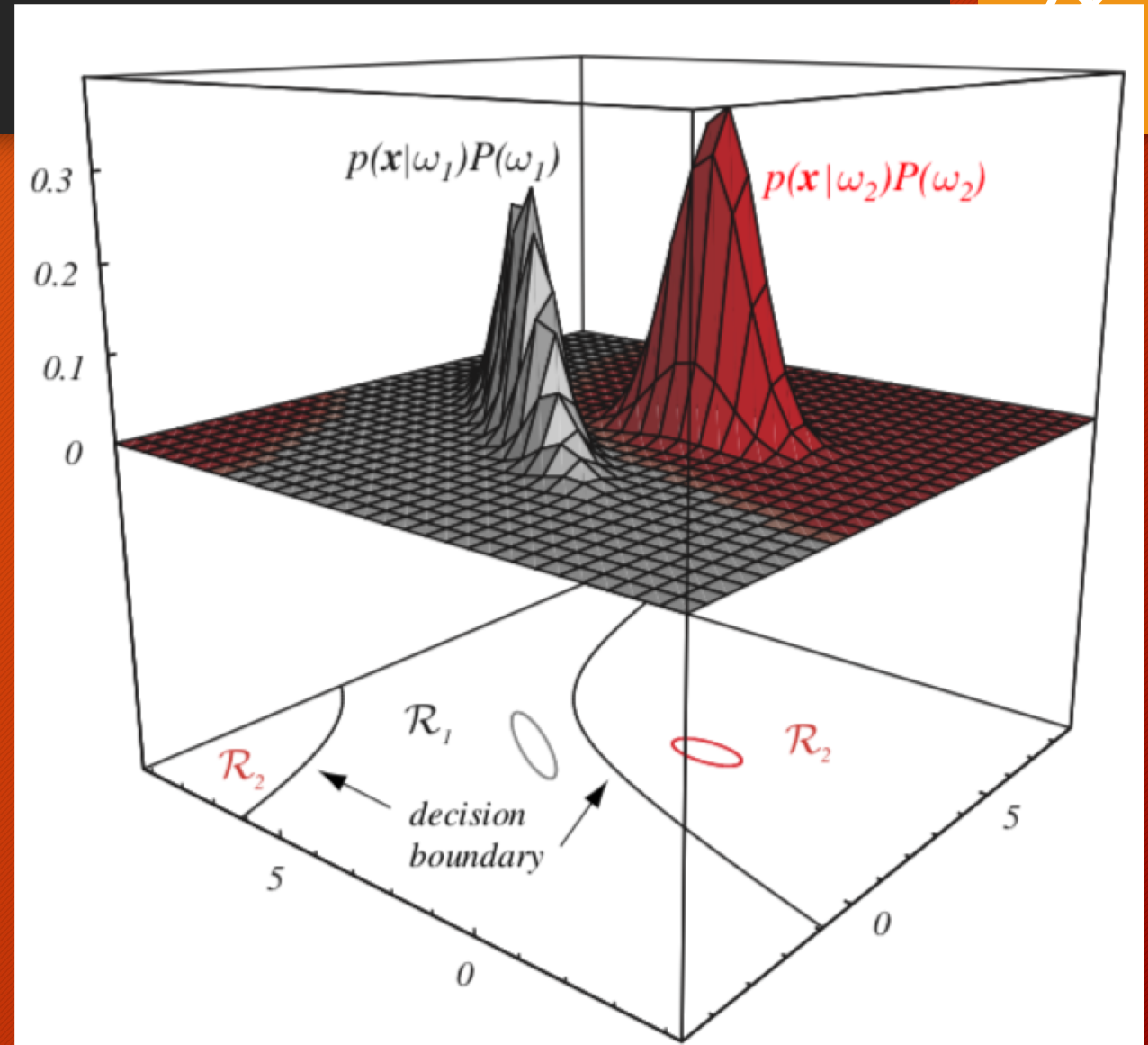
$$\Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma_2 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 2 \end{bmatrix}$$



## 4. 2D Mixture Model

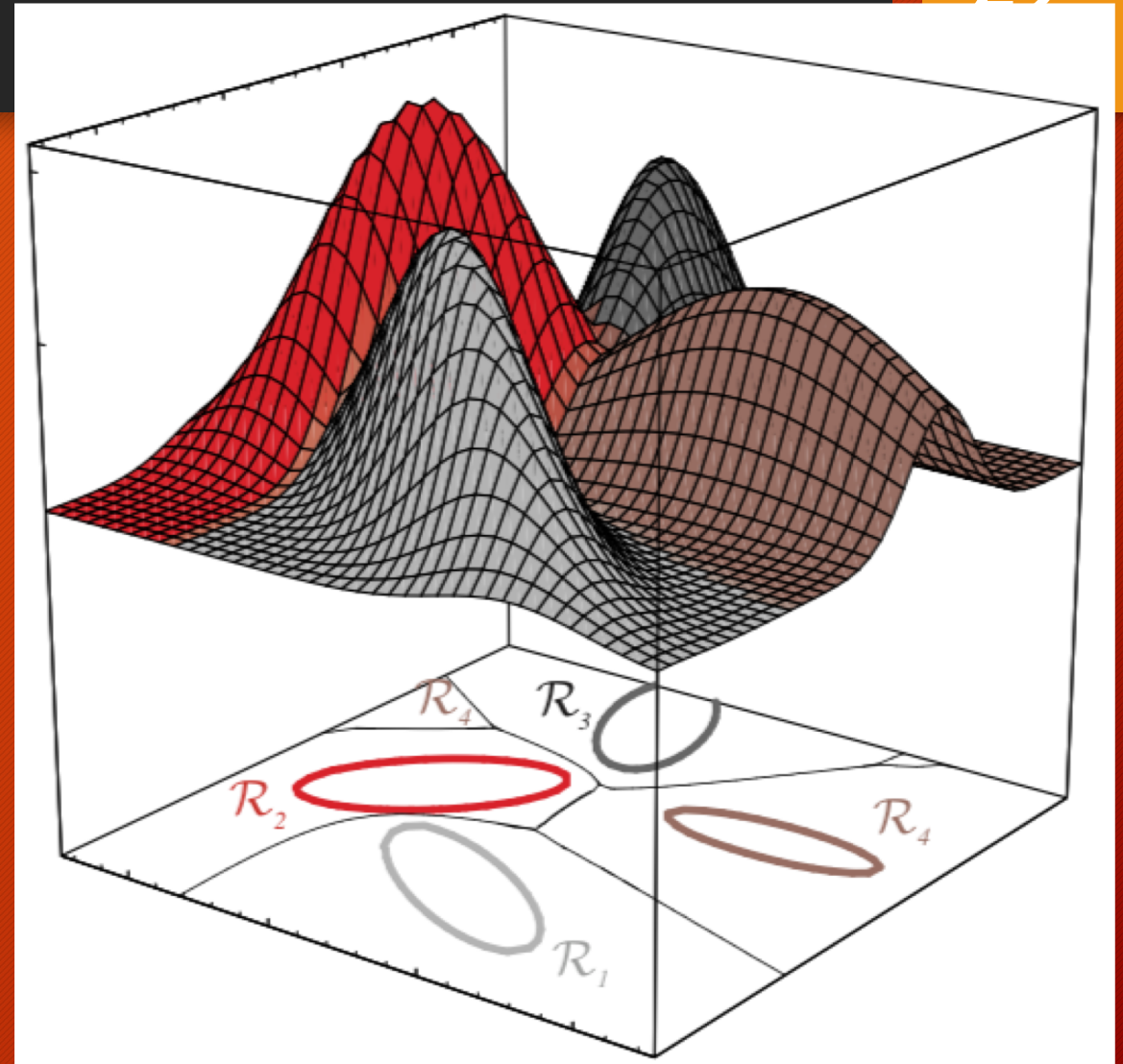
Decision boundary can be complex!!!  
(i.e. not always linear)



## 4. 2D Mixture Model

29

Multi-class decision  
boundaries



# 4. 2D Mixture Model

30

EM-based re-estimation formulas for parameters of a 2D MM are essentially the same:

$$\pi_k^{(\ell+1)} = \frac{1}{N} \sum_{n=1}^N \gamma_n^{(\ell)}(k)$$

$$\boldsymbol{\mu}_k^{(\ell+1)} = \frac{\sum_{n=1}^N \gamma_n^{(\ell)}(k) \mathbf{y}_n}{\sum_{n=1}^N \gamma_n^{(\ell)}(k)}$$

$$\boldsymbol{\Sigma}_k^{(\ell+1)} = \frac{\sum_{n=1}^N \gamma_n^{(\ell)}(k) \left( \mathbf{y}_n - \boldsymbol{\mu}_k^{(\ell+1)} \right)^T \left( \mathbf{y}_n - \boldsymbol{\mu}_k^{(\ell+1)} \right)}{\sum_{n=1}^N \gamma_n^{(\ell)}(k)}$$

- *Theory and Use of the EM Algorithm* By Maya R. Gupta and Yihua Chen, [Book pdf](#)
- *The EM algorithm and related statistical models* By Michiko Watanabe and Kazunori Yamaguchi, [Book pdf](#)
- *Pattern classification* by Richard O. Duda, Peter E. Hart and David G. Stork, 2015, [Book pdf](#), [Slides of the book](#)
- Finite Mixture model By Geoffrey McLachlan and D. Peel, [Book pdf](#)
- Python library for MM : [Pymix](#), [sklearn.mixture](#)

# Not seen!

32

- Variations about EM
  - GEM, CEM, SEM -- On-line EM , by O. Cappé
- Mixture of non-gaussian type:
  - M. of generalized hyperbolic distribution
  - M. of skew-normal distribution, M. of t-distribution
- Choosing the number of clusters via model selection criteria
  - BIC: Bayesian Information Criterion,
  - AIK: Akaike Information Criterion,
  - ICL: Integrated completed likelihood criterion

Biernacki, C., G. Celeux, and G. Govaert (2000). "Assessing a mixture model for clustering with the integrated completed likelihood". IEEe Trans. PAMI, Vol 22(7), pp. 719-725.