

Back to LDA and QDA

SY19 – Machine Learning

Chapter 7: Gaussian mixture models and the EM algorithm

Thierry Denœux

Université de technologie de Compiègne

<https://www.hds.utc.fr/~tdenoieux>
email: tdenoeux@utc.fr

Automne 2024



- In LDA and QDA, we assume that the conditional density of input vector X given $Y = k$ is multivariate Gaussian

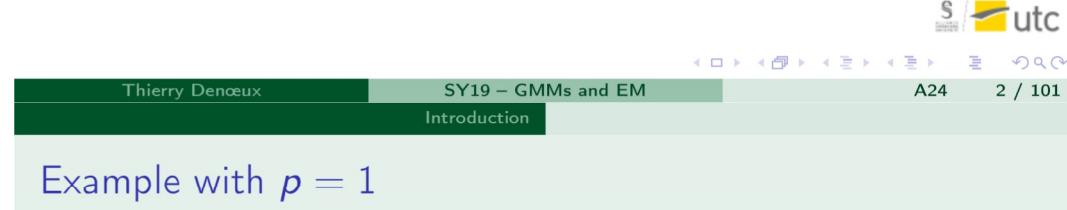
$$\phi(x; \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right)$$

(with $\Sigma_k = \Sigma$ in the case of LDA)

- The marginal density of X is then a mixture of c Gaussian densities:

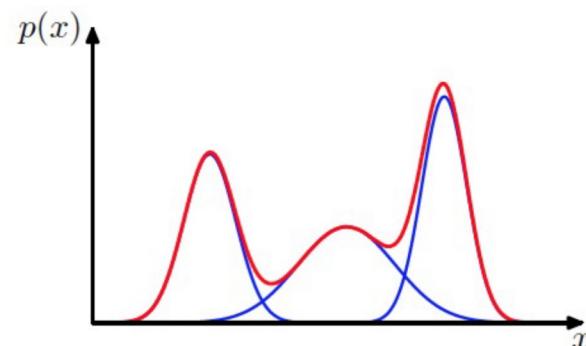
$$p(x) = \sum_{k=1}^c p(x \mid Y = k) P(Y = k) = \sum_{k=1}^c \pi_k \phi(x; \mu_k, \Sigma_k)$$

- This is called a Gaussian Mixture Model (GMM).

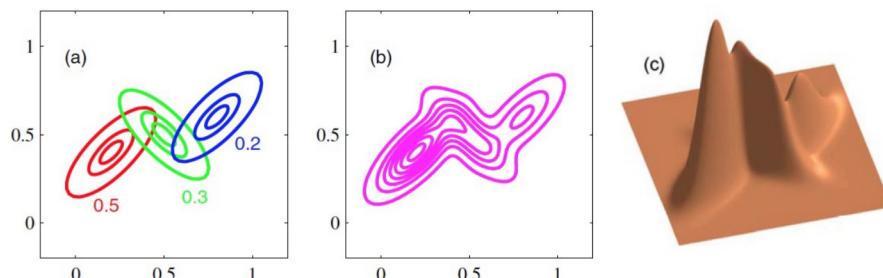


GMMs are widely used in Machine Learning for

- Density estimation
 - Clustering (finding groups in data)
 - Classification (modeling complex-shaped class distributions)
 - Regression (accounting for different linear relations within subgroups of a population)
 - etc.



Example with $p = 2$



Thierry Denœux SY19 – GMMs and EM A24 5 / 101

Introduction

Supervised learning

- In discriminant analysis, we observe both the input vector X and the response (class label) Y for n individuals taken randomly from the population.
- The learning set has the form $\mathcal{L}_s = \{(x_i, y_i)\}_{i=1}^n$. We say that the data are **labeled**.
- Learning a classifier from such data is called **supervised learning**.

How to generate data from a mixture?

- Assume $X \sim \sum_{k=1}^c \pi_k \mathcal{N}(\mu_k, \Sigma_k)$
- It is the marginal distribution of X in the pair (X, Y) , where
 - Y takes values in $\{1, \dots, c\}$ with probabilities π_1, \dots, π_c
 - The conditional distribution of X given $Y = k$ is the normal distribution $\mathcal{N}(\mu_k, \Sigma_k)$
- How to generate X ?
 - Generate $Y \in \{1, \dots, c\}$ with probabilities π_1, \dots, π_c
 - If $Y = k$, generate X from $p(x | Y = k) = \phi(x; \mu_k, \Sigma_k)$
- Remark: we can define mixtures of any distributions. In this chapter, we will focus (without loss of generality) on mixtures of normal distributions, called **Gaussian mixtures**.

Thierry Denœux SY19 – GMMs and EM A24 6 / 101

Introduction

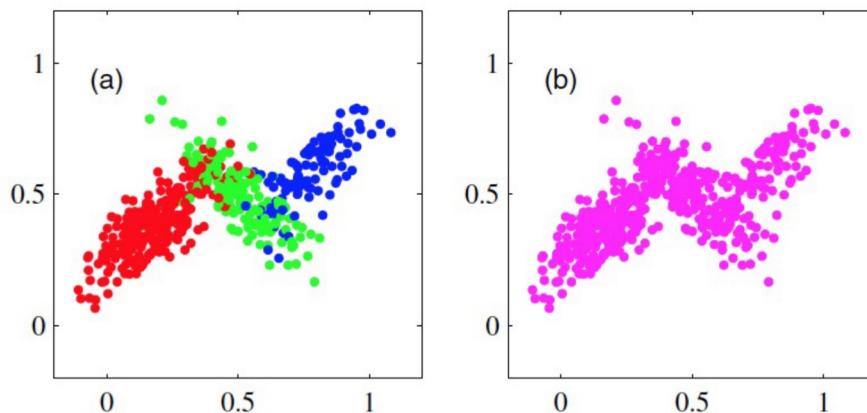
Unsupervised learning

- In some situations, we observe X , but Y is not observed. We say that Y is a **latent variable**.
- The learning set is composed of **unlabeled** data of the form $\mathcal{L}_{ns} = \{x_i\}_{i=1}^n$.
- Estimating the model parameters from such data is called **unsupervised learning**.
- Applications: density estimation, clustering, feature extraction.
- Unsupervised learning is usually more difficult than supervised learning, because we have less information to estimate the parameters.

Thierry Denœux SY19 – GMMs and EM A24 7 / 101

Thierry Denœux SY19 – GMMs and EM A24 8 / 101

Labeled vs. unlabeled data



Maximum likelihood: supervised case I

- In the case of **supervised learning** of GMMs, we have seen that the MLEs of μ_k , Σ_k and π_k have simple closed-form expressions.
- Assuming the sample $(X_1, Y_1), \dots, (X_n, Y_n)$ to be i.i.d., the likelihood function is

$$\begin{aligned} L(\theta; \mathcal{L}_s) &= \prod_{i=1}^n p(x_i, y_i) = \prod_{i=1}^n \underbrace{\frac{p(x_i | Y_i = y_i)}{\prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}}}}_{\text{term } \ell_k \text{ depending on } \mu_k \text{ and } \Sigma_k} \underbrace{P(Y_i = y_i)}_{\prod_{k=1}^c \pi_k^{y_{ik}}} \\ &= \prod_{i=1}^n \prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}} \pi_k^{y_{ik}} \end{aligned}$$

with $y_{ik} = I(y_i = k)$.

Semi-supervised learning

- Sometimes, we collect a lot of data, but we can label only a part of them.
- Examples: image data from the web, sensors on a robot, etc.
- The data then have the form

$$\mathcal{L}_{ss} = \underbrace{\{(x_i, y_i)\}_{i=1}^{n_s}}_{\text{labeled part}} \cup \underbrace{\{x_i\}_{i=n_s+1}^n}_{\text{unlabeled part}}$$

- Learning from such data is called **semi-supervised learning**.
- Semi-supervised learning is intermediate between supervised and unsupervised learning.

Maximum likelihood: supervised case II

- The log-likelihood function is

$$\ell(\theta; \mathcal{L}_s) = \sum_{k=1}^c \left\{ \underbrace{\sum_{i=1}^n y_{ik} \log \phi(x_i; \mu_k, \Sigma_k)}_{\text{term } \ell_k \text{ depending on } \mu_k \text{ and } \Sigma_k} \right\} + \underbrace{\sum_{i=1}^n \sum_{k=1}^c y_{ik} \log \pi_k}_{\text{term depending on } \pi_1, \dots, \pi_c}$$

- The parameters μ_k, Σ_k can be estimated separately using the data from class k .

MLE in the supervised case I

- We have

$$\ell_k = -\frac{1}{2} \sum_{i=1}^n y_{ik} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) - \frac{n_k}{2} \log |\Sigma_k| - \frac{n_k p}{2} \log(2\pi)$$

with $n_k = \sum_{i=1}^n y_{ik}$.

- The derivative wrt to μ_k is

$$\sum_i y_{ik} \Sigma_k^{-1} (x_i - \mu_k) = \Sigma_k^{-1} \sum_i y_{ik} (x_i - \mu_k).$$

Hence,

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^n y_{ik} x_i \quad \text{with} \quad n_k = \sum_{i=1}^n y_{ik}$$



Maximum likelihood: unsupervised case

- In the case of **unsupervised learning**, assuming the sample X_1, \dots, X_n to be i.i.d., the likelihood is

$$L(\theta; \mathcal{L}_{ns}) = \prod_{i=1}^n p(x_i)$$

and the log-likelihood function is

$$\ell(\theta; \mathcal{L}_{ns}) = \sum_{i=1}^n \log p(x_i) = \sum_{i=1}^n \left(\log \sum_{k=1}^c \pi_k \phi(x_i; \mu_k, \Sigma_k) \right)$$

- We can no longer separate the terms corresponding to each class.
- Maximizing the log-likelihood becomes a difficult nonlinear optimization problem, for which no closed-form solution exists.
- A powerful method: the **Expectation-Maximization (EM)** algorithm.



MLE in the supervised case II

- It can be shown that

$$\widehat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^n y_{ik} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

- To find the MLE of the π_k , we maximize the last term

$$\sum_{i=1}^n \sum_{k=1}^c y_{ik} \log \pi_k$$

wrt to π_k , subject to the constraint $\sum_{k=1}^c \pi_k = 1$.

- The solution is

$$\widehat{\pi}_k = \frac{n_k}{n}, \quad k = 1, \dots, c$$



Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts



EM Algorithm

- An iterative optimization strategy useful when maximizing the likelihood is difficult, but:
 - ▶ There are **missing** (non-observed) data
 - ▶ If the missing data were observed, maximizing the likelihood would be easy.
- Many applications in statistics and ML.
- Can be very simple to implement. Can reliably find an optimum through stable, uphill steps.

Thierry Denœux SY19 – GMMs and EM A24 17 / 101

EM algorithm General formulation

Notation

- \mathbf{X} : Observed variables
- \mathbf{Y} : Missing or latent variables
- \mathbf{Z} : Complete data $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$
- θ : Unknown parameter
- $L(\theta)$: observed-data likelihood, short for $L(\theta; \mathbf{x}) = p(\mathbf{x}; \theta)$
- $L_c(\theta)$: complete-data likelihood, short for $L(\theta; \mathbf{z}) = p(\mathbf{z}; \theta)$
- $\ell(\theta), \ell_c(\theta)$: observed and complete-data log-likelihoods

Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts

Thierry Denœux SY19 – GMMs and EM A24 18 / 101

EM algorithm General formulation

Q function

- Suppose we seek to maximize $L(\theta)$ with respect to θ .
- Define $Q(\theta; \theta^{(t)})$ to be the **expectation of the complete-data log-likelihood** (assuming $\theta = \theta^{(t)}$), conditional on the observed data $\mathbf{X} = \mathbf{x}$. Namely

$$\begin{aligned} Q(\theta, \theta^{(t)}) &= \mathbb{E}_{\theta^{(t)}} \{ \ell_c(\theta) | \mathbf{x} \} \\ &= \mathbb{E}_{\theta^{(t)}} \{ \log p(\mathbf{Z}; \theta) | \mathbf{x} \} \\ &= \int [\log p(\mathbf{z}; \theta)] \underbrace{p(\mathbf{z} | \mathbf{x}; \theta^{(t)})}_{p(\mathbf{y} | \mathbf{x}; \theta^{(t)})} d\mathbf{y} \end{aligned}$$

$(p(\mathbf{z} | \mathbf{x}; \theta^{(t)}) = p(\mathbf{y} | \mathbf{x}; \theta^{(t)})$ because \mathbf{Y} is the only random part of \mathbf{Z} once we are given $\mathbf{X} = \mathbf{x}$)

Thierry Denœux SY19 – GMMs and EM A24 19 / 101

Thierry Denœux SY19 – GMMs and EM A24 20 / 101

The EM Algorithm

Start with $\theta^{(0)}$ and set $t = 0$. Then

- ① **E step:** Compute $Q(\theta, \theta^{(t)})$.
- ② **M step:** Maximize $Q(\theta, \theta^{(t)})$ with respect to θ ; set $\theta^{(t+1)}$ equal to the maximizer of Q .
- ③ Return to the E step and increment t unless a stopping criterion has been met, e.g.,

$$|\ell(\theta^{(t+1)}) - \ell(\theta^{(t)})| \leq \epsilon$$

or

$$\|\theta^{(t+1)} - \theta^{(t)}\| \leq \epsilon$$

Thierry Denœux SY19 – GMMs and EM A24 21 / 101

EM algorithm Simple example

Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts

Convergence of the EM Algorithm

- We will show that $L(\theta)$ increases after each EM iteration, i.e., $L(\theta^{(t+1)}) \geq L(\theta^{(t)})$ for $t = 0, 1, \dots$ (see below)
- Consequently, the algorithm converges to a **local maximum** of $L(\theta)$ if the likelihood function is bounded above.
- Typically, we run the algorithm several times with random initial conditions, and we keep the results of the best run.

Thierry Denœux SY19 – GMMs and EM A24 22 / 101

EM algorithm Simple example

Mixture of two univariate normal distributions

- Let $\mathbf{X} = (X_1, \dots, X_n)$ be an i.i.d. sample from a mixture of two univariate normal distributions $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$, with pdf

$$p(x_i; \theta) = \pi\phi(x_i; \mu_1, \sigma_1^2) + (1 - \pi)\phi(x_i; \mu_2, \sigma_2^2)$$

where $\phi(\cdot; \mu, \sigma^2)$ is the univariate normal pdf and

$$\theta = (\mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \pi)^T$$

is the vector of parameters.

- We introduce **latent variables** $\mathbf{Y} = (Y_1, \dots, Y_n)$, such that

- ▶ $Y_i \sim \mathcal{B}(\pi)$
- ▶ $p(x_i | Y_i = 1) = \phi(x_i; \mu_1, \sigma_1^2)$ and
- ▶ $p(x_i | Y_i = 0) = \phi(x_i; \mu_2, \sigma_2^2)$.

Thierry Denœux SY19 – GMMs and EM A24 23 / 101

Thierry Denœux SY19 – GMMs and EM A24 24 / 101

Observed and complete-data likelihoods

- Observed-data likelihood:

$$L(\theta) = \prod_{i=1}^n p(x_i; \theta) = \prod_{i=1}^n [\pi\phi(x_i; \mu_1, \sigma_1^2) + (1 - \pi)\phi(x_i; \mu_2, \sigma_2^2)]$$

- Complete-data likelihood:

$$\begin{aligned} L_c(\theta) &= \prod_{i=1}^n p(x_i, y_i; \theta) = \prod_{i=1}^n p(x_i | y_i; \theta)p(y_i; \pi) \\ &= \prod_{i=1}^n \{\phi(x_i; \mu_1, \sigma_1^2)^{y_i} \phi(x_i; \mu_2, \sigma_2^2)^{1-y_i} \pi^{y_i} (1-\pi)^{1-y_i}\} \end{aligned}$$

Thierry Denœux SY19 – GMMs and EM A24 25 / 101

EM algorithm: E-step

Compute

$$\begin{aligned} y_i^{(t)} &= \mathbb{E}_{\theta^{(t)}}[Y_i | x_i] \\ &= \mathbb{P}_{\theta^{(t)}}[Y_i = 1 | x_i] \\ &= \frac{\phi(x_i; \mu_1^{(t)}, (\sigma_1^2)^{(t)})\pi^{(t)}}{\phi(x_i; \mu_1^{(t)}, (\sigma_1^2)^{(t)})\pi^{(t)} + \phi(x_i; \mu_2^{(t)}, (\sigma_2^2)^{(t)})(1 - \pi^{(t)})} \end{aligned}$$

Thierry Denœux SY19 – GMMs and EM A24 27 / 101

Derivation of function Q

- Complete-data log-likelihood:

$$\begin{aligned} \ell_c(\theta) &= \sum_{i=1}^n \{y_i \log \phi(x_i; \mu_1, \sigma_1^2) + (1 - y_i) \log \phi(x_i; \mu_2, \sigma_2^2)\} \\ &\quad + \sum_{i=1}^n \{y_i \log \pi + (1 - y_i) \log(1 - \pi)\} \end{aligned}$$

- It is linear in the y_i . Consequently, the Q function is simply

$$\begin{aligned} Q(\theta, \theta^{(t)}) &= \sum_{i=1}^n \left\{ y_i^{(t)} \log \phi(x_i; \mu_1, \sigma_1^2) + (1 - y_i^{(t)}) \log \phi(x_i; \mu_2, \sigma_2^2) \right\} \\ &\quad + \sum_{i=1}^n \left\{ y_i^{(t)} \log \pi + (1 - y_i^{(t)}) \log(1 - \pi) \right\} \end{aligned}$$

with $y_i^{(t)} = \mathbb{E}_{\theta^{(t)}}[Y_i | x_i]$.

Thierry Denœux SY19 – GMMs and EM A24 26 / 101

EM algorithm: M-step

Maximize $Q(\theta, \theta^{(t)})$. We get

$$\begin{aligned} \pi^{(t+1)} &= \frac{n_1^{(t)}}{n} \\ \mu_1^{(t+1)} &= \frac{\sum_{i=1}^n y_i^{(t)} x_i}{n_1^{(t)}}, (\sigma_1^2)^{(t+1)} = \frac{\sum_{i=1}^n y_i^{(t)} (x_i - \mu_1^{(t+1)})^2}{n_1^{(t)}} \\ \mu_2^{(t+1)} &= \frac{\sum_{i=1}^n (1 - y_i^{(t)}) x_i}{n_2^{(t)}}, (\sigma_2^2)^{(t+1)} = \frac{\sum_{i=1}^n (1 - y_i^{(t)}) (x_i - \mu_2^{(t+1)})^2}{n_2^{(t)}} \end{aligned}$$

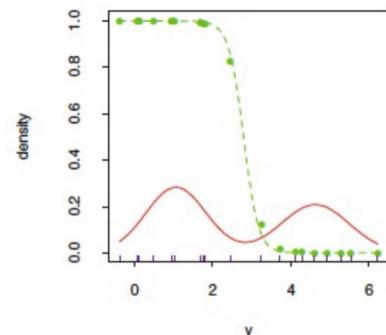
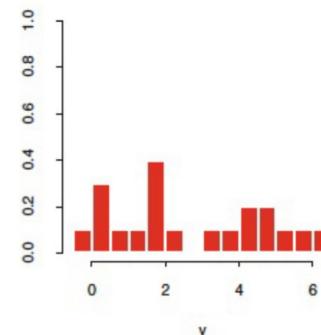
with

$$n_1^{(t)} = \sum_{i=1}^n y_i^{(t)} \quad \text{and} \quad n_2^{(t)} = n - n_1^{(t)}$$

Thierry Denœux SY19 – GMMs and EM A24 28 / 101

Example

-0.39	0.12	0.94	1.67	1.76	2.44	3.72	4.28	4.92	5.53
0.06	0.48	1.01	1.68	1.80	3.25	4.12	4.60	5.28	6.22



(green curve: $\hat{P}_\theta[Y = 1 | x]$ as a function of x , assuming $Y = 1$ corresponds to the left component)

Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

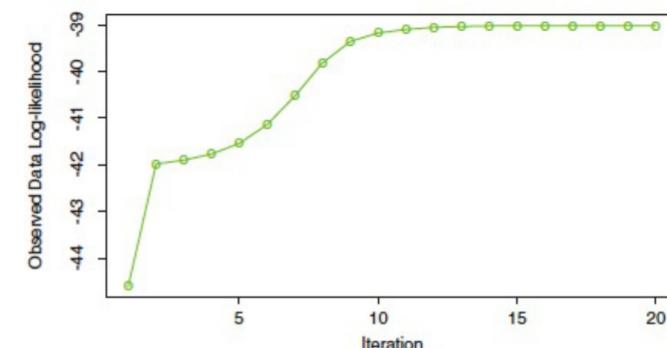
2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts

Example (continued)



Solution: $\hat{\mu}_1 = 4.66$, $\hat{\sigma}_1 = 0.91$, $\hat{\mu}_2 = 1.08$, $\hat{\sigma}_2 = 0.90$, $\hat{\pi} = 0.45$.

Why does it work?

1 Optimization transfer:

- ▶ At each step t , EM transfers optimization from ℓ to a surrogate function $G(\cdot, \theta^{(t)})$ defined as

$$G(\theta, \theta^{(t)}) = Q(\theta, \theta^{(t)}) + \ell(\theta^{(t)}) - Q(\theta^{(t)}, \theta^{(t)})$$

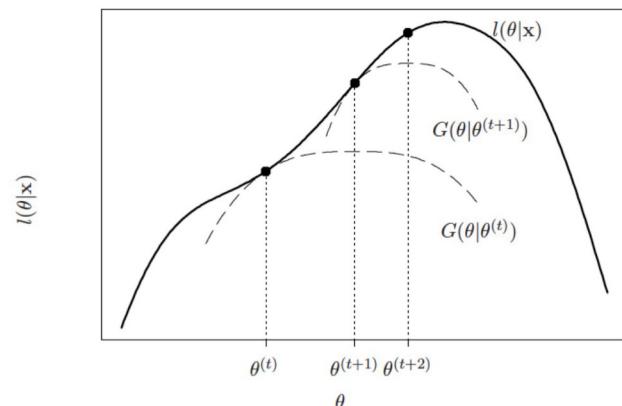
(Remark: The last two terms in $G(\theta, \theta^{(t)})$ do not depend on θ , so Q and G are maximized at the same θ).

- ▶ Each function $G(\cdot, \theta^{(t)})$ has the following properties:
 - It lies everywhere below ℓ : $\forall \theta$, $\ell(\theta) \geq G(\theta, \theta^{(t)})$
 - It is tangent to ℓ at $\theta^{(t)}$.

We say that G is a **minorizing function** for ℓ (see next slide).

- 2 As a result, each M-step increases the log-likelihood.

The nature of EM (continued)



One-dimensional illustration of EM algorithm as a minorization or optimization transfer strategy. Each E step forms a minorizing function G , and each M step maximizes it to provide an uphill step.

Proof

- We have

$$p(y | x; \theta) = \frac{p(x, y; \theta)}{p(x; \theta)} = \frac{p(z; \theta)}{p(x; \theta)} \Rightarrow p(x; \theta) = \frac{p(z; \theta)}{p(y | x; \theta)}$$

- Consequently,

$$\ell(\theta) = \log p(x; \theta) = \underbrace{\log p(z; \theta)}_{\ell_c(\theta)} - \log p(y | x; \theta)$$

- Taking expectations on both sides wrt the conditional distribution of Z given $X = x$ and using $\theta^{(t)}$ for θ :

$$\ell(\theta) = Q(\theta, \theta^{(t)}) - \underbrace{\mathbb{E}_{\theta^{(t)}}[\log p(Y | x; \theta) | x]}_{H(\theta, \theta^{(t)})} \quad (1)$$

Proof: $\ell(\cdot)$ dominates $G(\cdot, \theta^{(t)})$

- Now, for all $\theta \in \Theta$,

$$H(\theta, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}) = \mathbb{E}_{\theta^{(t)}} \left[\log \frac{p(Y | x; \theta)}{p(Y | x; \theta^{(t)})} \mid x \right] \quad (2a)$$

$$\leq \log \mathbb{E}_{\theta^{(t)}} \left[\underbrace{\frac{p(Y | x; \theta)}{p(Y | x; \theta^{(t)})}}_{\int \frac{p(y | x; \theta)}{p(y | x; \theta^{(t)})} p(y | x; \theta^{(t)}) dy} \mid x \right] (*) \quad (2b)$$

$$\leq \log \underbrace{\int p(y | x; \theta) dy}_{1} = 0 \quad (2c)$$

Hence, for all $\theta \in \Theta$,

$$H(\theta^{(t)}, \theta^{(t)}) \geq H(\theta, \theta^{(t)})$$

$$Q(\theta^{(t)}, \theta^{(t)}) - \ell(\theta^{(t)}) \geq Q(\theta, \theta^{(t)}) - \ell(\theta)$$

$$\ell(\theta) \geq \underbrace{Q(\theta, \theta^{(t)}) + \ell(\theta^{(t)}) - Q(\theta^{(t)}, \theta^{(t)})}_{G(\theta, \theta^{(t)})}$$

(*): from the concavity of the log and Jensen's inequality.

- Hence, $\theta^{(t)}$ is a maximizer of $H(\theta, \theta^{(t)})$

Proof: G is tangent to ℓ at $\theta^{(t)}$

- As $\theta^{(t)}$ maximizes $H(\theta, \theta^{(t)}) = Q(\theta, \theta^{(t)}) - \ell(\theta)$, we have

$$H'(\theta, \theta^{(t)})|_{\theta=\theta^{(t)}} = Q'(\theta, \theta^{(t)})|_{\theta=\theta^{(t)}} - \ell'(\theta)|_{\theta=\theta^{(t)}} = 0$$

so

$$Q'(\theta, \theta^{(t)})|_{\theta=\theta^{(t)}} = \ell'(\theta)|_{\theta=\theta^{(t)}}$$

- Consequently, as $G(\theta, \theta^{(t)}) = Q(\theta, \theta^{(t)}) + \text{cst}$,

$$G'(\theta, \theta^{(t)})|_{\theta=\theta^{(t)}} = Q'(\theta, \theta^{(t)})|_{\theta=\theta^{(t)}} = \ell'(\theta)|_{\theta=\theta^{(t)}}$$

The screenshot shows a presentation slide with the following details:

- Header:** Thierry Denœux, SY19 – GMMs and EM, Parameter estimation in GMMs.
- Navigation:** A toolbar with various icons for navigating through the presentation.
- Page Information:** A24, 37 / 101.
- Content:** The slide title is "Overview".

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts



Proof: monotonicity

- From (1),

$$\ell(\theta^{(t+1)}) - \ell(\theta^{(t)}) = \underbrace{Q(\theta^{(t+1)}, \theta^{(t)}) - Q(\theta^{(t)}, \theta^{(t)})}_A - \underbrace{(H(\theta^{(t+1)}, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}))}_B$$

- $A \geq 0$ because $\theta^{(t+1)}$ is a maximizer of $Q(\theta, \theta^{(t)})$, and $B \leq 0$ because from (2) $\theta^{(t)}$ is a maximizer of $H(\theta, \theta^{(t)})$.

- Hence,

$$\ell(\theta^{(t+1)}) \geq \ell(\theta^{(t)})$$

The screenshot shows a presentation slide with the following details:

- Header:** Thierry Denœux, SY19 – GMMs and EM, Parameter estimation in GMMs | Unsupervised learning.
- Navigation:** A toolbar with various icons for navigating through the presentation.
- Page Information:** A24, 38 / 101.
- Content:** The slide title is "Overview".

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

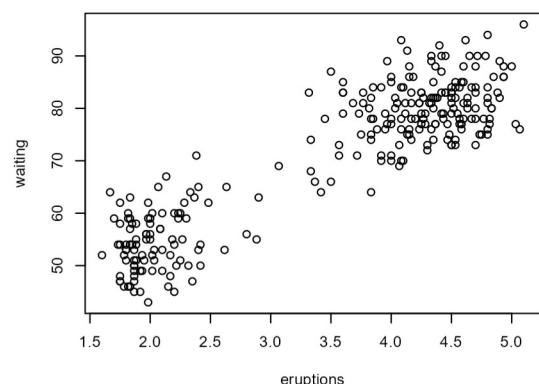
- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts



Old Faithful geyser data



Waiting time between eruptions and duration of the eruption (in min) for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA (272 observations).

General GMM

- Let $\mathbf{X} = (X_1, \dots, X_n)$ be an i.i.d. sample from a mixture of c multivariate normal distributions $\mathcal{N}(\mu_k, \Sigma_k)$ with proportions π_k . The pdf of X_i is

$$p(x_i; \theta) = \sum_{k=1}^c \pi_k \phi(x_i; \mu_k, \Sigma_k)$$

where θ is the vector of parameters.

- We introduce latent variables $\mathbf{Y} = (Y_1, \dots, Y_n)$, such that
 - Y_i has a categorical distribution with parameters (π_1, \dots, π_c)
 - $p(x_i | Y_i = k) = \phi(x_i; \mu_k, \Sigma_k)$, $k = 1 \dots, c$

Old Faithful geyser data (continued)

- Questions:

- How can we best partition these data into c groups/clusters (for instance, $c = 2$)?
- What is the most plausible number of groups?

- Approach:

- Fit GMMs to these data
- Select the best model according to some criterion

Observed and complete-data likelihoods

- Observed-data likelihood:

$$L(\theta) = \prod_{i=1}^n p(x_i; \theta) = \prod_{i=1}^n \sum_{k=1}^c \pi_k \phi(x_i; \mu_k, \Sigma_k)$$

- Complete-data likelihood:

$$\begin{aligned} L_c(\theta) &= \prod_{i=1}^n p(x_i, y_i; \theta) = \prod_{i=1}^n p(x_i | y_i; \theta) p(y_i; \pi) \\ &= \prod_{i=1}^n \prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}} \pi_k^{y_{ik}} \end{aligned}$$

Derivation of function Q

- Complete-data log-likelihood:

$$\ell_c(\theta) = \sum_{i=1}^n \sum_{k=1}^c y_{ik} \log \phi(x_i; \mu_k, \Sigma_k) + \sum_{i=1}^n \sum_{k=1}^c y_{ik} \log \pi_k$$

- It is linear in the y_{ik} . Consequently, the Q function is simply

$$Q(\theta, \theta^{(t)}) = \underbrace{\sum_{k=1}^c \sum_{i=1}^n y_{ik}^{(t)} \log \phi(x_i; \mu_k, \Sigma_k)}_{\text{term depending only on } \mu_k \text{ and } \Sigma_k} + \underbrace{\sum_{i=1}^n \sum_{k=1}^c y_{ik}^{(t)} \log \pi_k}_{\text{term depending only on } \{\pi_k\}}$$

with $y_{ik}^{(t)} = \mathbb{E}_{\theta^{(t)}}[Y_{ik} | x_i] = \mathbb{P}_{\theta^{(t)}}[Y_i = k | x_i]$.



GMM with the package mclust

```
library(mclust)
data(faithful)

faithfulMclust <- Mclust(faithful, G=2, modelNames="VVV")
plot(faithfulMclust)
```



EM algorithm

- E-step: compute

$$\begin{aligned} y_{ik}^{(t)} &= \mathbb{P}_{\theta^{(t)}}[Y_i = k | x_i] \\ &= \frac{\phi(x_i; \mu_k^{(t)}, \Sigma_k^{(t)}) \pi_k^{(t)}}{\sum_{\ell=1}^c \phi(x_i; \mu_\ell^{(t)}, \Sigma_\ell^{(t)}) \pi_\ell^{(t)}} \end{aligned}$$

- M-step: Maximize $Q(\theta, \theta^{(t)})$. We get

$$\pi_k^{(t+1)} = \frac{n_k^{(t)}}{n}, \quad \mu_k^{(t+1)} = \frac{1}{n_k^{(t)}} \sum_{i=1}^n y_{ik}^{(t)} x_i$$

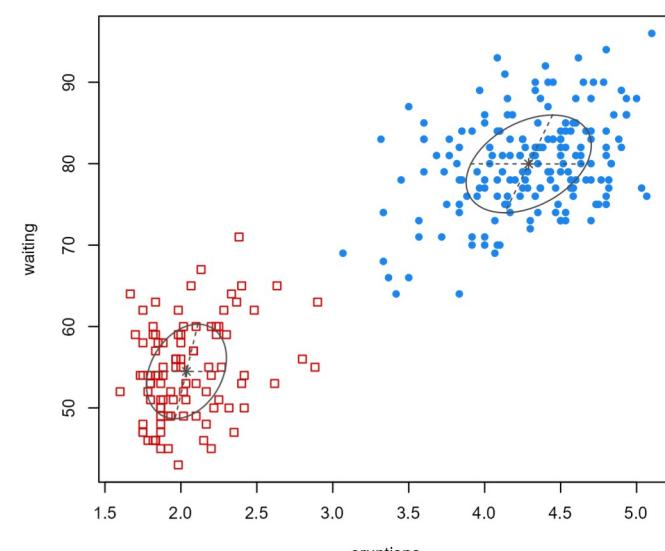
$$\Sigma_k^{(t+1)} = \frac{1}{n_k^{(t)}} \sum_{i=1}^n y_{ik}^{(t)} (x_i - \mu_k^{(t+1)}) (x_i - \mu_k^{(t+1)})^T$$

with $n_k^{(t)} = \sum_{i=1}^n y_{ik}^{(t)}$.



Result

Classification



Choosing the number of clusters

- In clustering, selecting the number of clusters is often a difficult problem.
- With the GMM approach, this is a model selection problem. We can use the BIC criterion. (Reminder: $BIC = -2\ell(\hat{\theta}) + d \log(n)$; actually, Mclust computes $-BIC$).

```
> faithfulMclust <- Mclust(faithful, modelNames="VVV")
> summary(faithfulMclust)
```

Gaussian finite mixture model fitted by EM algorithm

Mclust VVV (ellipsoidal, varying volume, shape, and orientation) model with 2 components:

log.likelihood	n	df	BIC	ICL
-1130.264	272	11	-2322.192	-2322.695

Clustering table:

1	2
175	97

Thierry Denœux

SY19 – GMMs and EM

A24 49 / 101

Parameter estimation in GMMs

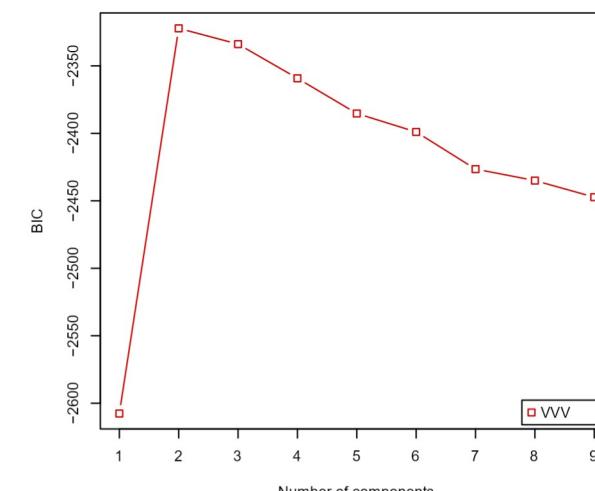
Unsupervised learning

Reducing the number of parameters

- The general model has $c[p + p(p + 1)/2 + 1] - 1$ parameters.
- When n is small and/or p is large: we need more parsimonious models (i.e., models with fewer parameters).
- Simple approaches:
 - Assume equal covariance matrix (homoscedasticity)
 - Assume the covariance matrices to be diagonal, or scalar
- More flexible approach: reparameterize matrix Σ_k using its eigendecomposition.

Choosing the number of clusters

`plot(faithfulMclust)`



Thierry Denœux

SY19 – GMMs and EM

A24 50 / 101

Parameter estimation in GMMs

Unsupervised learning

Eigendecomposition of Σ_k

We have seen that, as Σ_k is symmetric, it can be decomposed as

$$\Sigma_k = D_k \Lambda_k D_k^T$$

where

- $\Lambda_k = \text{diag}(\lambda_{k1}, \dots, \lambda_{kp})$ is a real diagonal matrix whose components are the **decreasing eigenvalues** of Σ_k , with $|\Lambda_k| = \prod_{j=1}^p \lambda_{kj} = |\Sigma_k|$
- D_k is an orthogonal matrix ($D_k D_k^T = I_p$) whose columns are the **normalized eigenvectors** of Σ_k ; it is a rotation matrix

Decomposition of Λ_k

Matrix Λ_k can be further decomposed as

$$\Lambda_k = \lambda_k \mathbf{A}_k$$

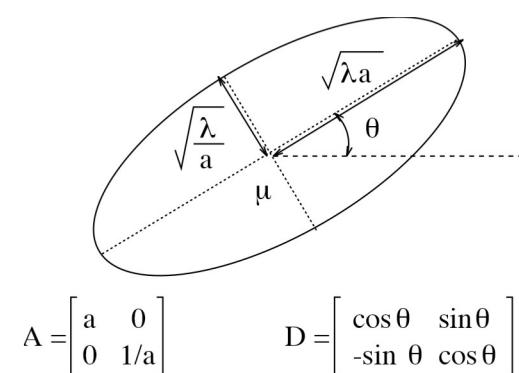
where

- $\lambda_k = \left(\prod_{j=1}^p \lambda_{kj} \right)^{1/p} = |\Sigma_k|^{1/p}$
- $\mathbf{A}_k = \Lambda_k / \lambda_k$ is a diagonal matrix verifying $|\mathbf{A}_k| = 1$.

Thierry Denœux SY19 – GMMs and EM A24 53 / 101

Parameter estimation in GMMs Unsupervised learning

Example in \mathbb{R}^2



- \mathbf{D} : rotation matrix, angle θ
- \mathbf{A} : diagonal matrix with diagonal terms a and $1/a$
- The eigenvalues of Σ are λa and λ/a .

Interpretation

- Each term in the decomposition

$$\Sigma_k = \lambda_k \mathbf{D}_k \mathbf{A}_k \mathbf{D}_k^T$$

has a simple interpretation:

- ▶ \mathbf{A}_k describes the **shape** of the cluster (defined by the ratios of the eigenvalues of Σ_k)
- ▶ \mathbf{D}_k (a rotation matrix) describes its **orientation**
- ▶ λ_k describes its **volume**

- Number of parameters:

Σ_k	λ_k	\mathbf{A}_k	\mathbf{D}_k
$\frac{p(p+1)}{2}$	1	$p - 1$	$\frac{p(p-1)}{2}$

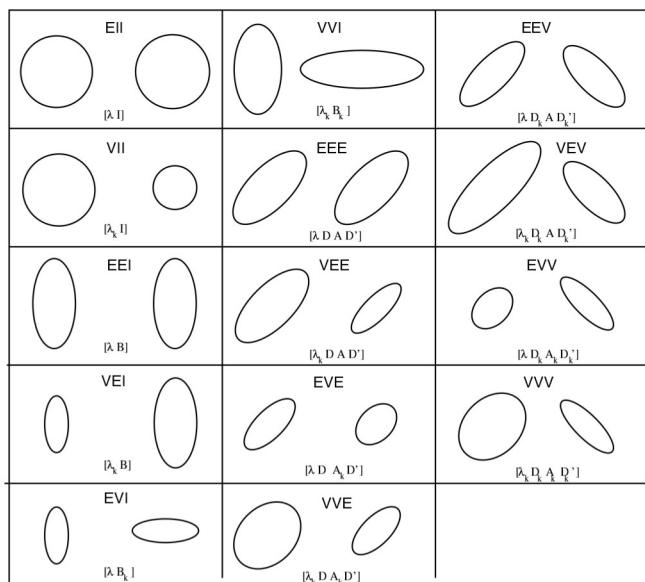
Thierry Denœux SY19 – GMMs and EM A24 54 / 101

Parameter estimation in GMMs Unsupervised learning

Parsimonious models

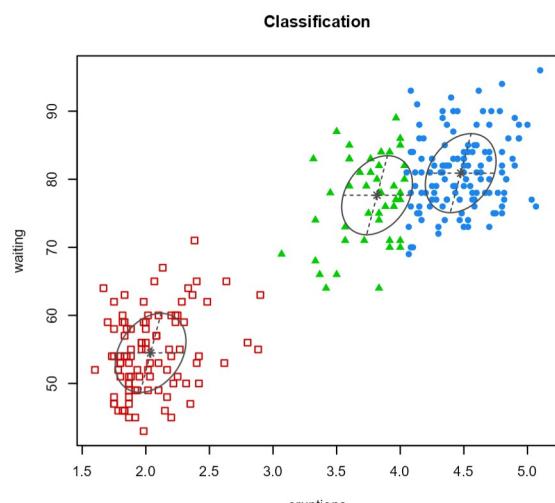
- With this parametrization, the parameters of the GMM are: the centers, volumes, shapes, orientations and proportions.
- 28 different models:
 - ▶ Spherical, diagonal, arbitrary
 - ▶ Volumes equal or not
 - ▶ Shapes equal or not
 - ▶ Orientations equal or not
 - ▶ Proportions equal or not

The 14 models based on assumptions on variance matrices



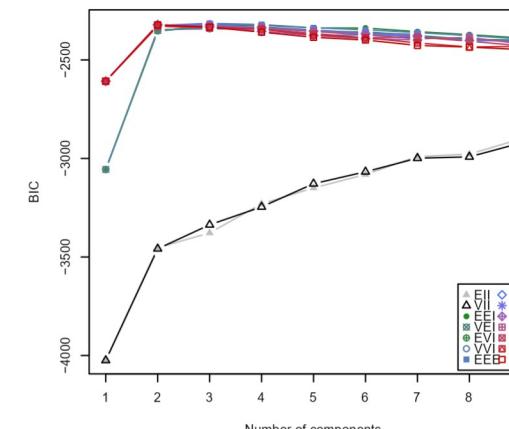
Best model

Best model: EEE or λDAD^T (ellipsoidal, equal volume, shape and orientation) model with 3 components



Parsimonious models in mclust

```
faithfulMcclus <- Mcclus(faithful)
plot(faithfulMcclus)
```



Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts

Model

- In semi-supervised learning, the data have the form

$$\mathcal{L}_{ss} = \underbrace{\{(x_i, y_i)\}_{i=1}^{n_s}}_{\text{labeled part}} \cup \underbrace{\{x_i\}_{i=n_s+1}^n}_{\text{unlabeled part}}$$

- Observed-data likelihood:

$$\begin{aligned} L(\theta) &= \prod_{i=1}^{n_s} p(x_i, y_i; \theta) \prod_{i=n_s+1}^n p(x_i; \theta) \\ &= \left(\prod_{i=1}^{n_s} \prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}} \pi_k^{y_{ik}} \right) \left(\prod_{i=n_s+1}^n \sum_{k=1}^c \pi_k \phi(x_i; \mu_k, \Sigma_k) \right) \end{aligned}$$

Q function I

- Complete-data log-likelihood:

$$\begin{aligned} \ell_c(\theta) &= \sum_{i=1}^{n_s} \sum_{k=1}^c y_{ik} (\log \phi(x_i; \mu_k, \Sigma_k) + \log \pi_k) + \\ &\quad \sum_{i=n_s+1}^n \sum_{k=1}^c y_{ik} (\phi(x_i; \mu_k, \Sigma_k) + \log \pi_k) \end{aligned}$$

Q function I

- Complete-data likelihood:

$$\begin{aligned} L_c(\theta) &= \prod_{i=1}^n \prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}} \pi_k^{y_{ik}} \\ &= \underbrace{\prod_{i=1}^{n_s} \prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}} \pi_k^{y_{ik}}}_{\text{observed}} \underbrace{\prod_{i=n_s+1}^n \prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}} \pi_k^{y_{ik}}}_{\text{non-observed}} \end{aligned}$$

Q function III

- Q function:

$$\begin{aligned} Q(\theta, \theta^{(t)}) &= \sum_{i=1}^{n_s} \sum_{k=1}^c y_{ik} (\log \phi(x_i; \mu_k, \Sigma_k) + \log \pi_k) + \\ &\quad \sum_{i=n_s+1}^n \sum_{k=1}^c y_{ik}^{(t)} (\log \phi(x_i; \mu_k, \Sigma_k) + \log \pi_k) \\ &= \sum_{k=1}^c \sum_{i=1}^n y_{ik}^{(t)} \log \phi(x_i; \mu_k, \Sigma_k) + \sum_{i=1}^n \sum_{k=1}^c y_{ik}^{(t)} \log \pi_k \end{aligned}$$

with

$$y_{ik}^{(t)} = \begin{cases} y_{ik} & i = 1, \dots, n_s \\ \mathbb{E}_{\theta^{(t)}}[Y_{ik} | x_i] & i = n_s + 1, \dots, n \end{cases}$$

EM algorithm

E-step: Compute

$$y_{ik}^{(t)} = \begin{cases} y_{ik} & i = 1, \dots, n_s \text{ (fixed)} \\ \frac{\phi(x_i; \mu_k^{(t)}, \Sigma_k^{(t)}) \pi_k^{(t)}}{\sum_{\ell=1}^c \phi(x_i; \mu_\ell^{(t)}, \Sigma_\ell^{(t)}) \pi_\ell^{(t)}} & i = n_s + 1, \dots, n \end{cases}$$

M-step: Same as in the unsupervised case.

$$\pi_k^{(t+1)} = \frac{n_k^{(t)}}{n}, \quad \mu_k^{(t+1)} = \frac{1}{n_k^{(t)}} \sum_{i=1}^n y_{ik}^{(t)} x_i$$

$$\Sigma_k^{(t+1)} = \frac{1}{n_k^{(t)}} \sum_{i=1}^n y_{ik}^{(t)} (x_i - \mu_k^{(t+1)}) (x_i - \mu_k^{(t+1)})^T$$

$$\text{with } n_k^{(t)} = \sum_{i=1}^n y_{ik}^{(t)}.$$



Mixture Discriminant Analysis

- GMM can also be useful in supervised classification.
- Here, we model the distribution of X in each class by a GMM:

$$p(x | Y = k) = \sum_{r=1}^{R_k} \pi_{kr} \phi(x; \mu_{kr}, \Sigma_{kr})$$

$$\text{with } \sum_{r=1}^{R_k} \pi_{kr} = 1.$$

- This method is called Mixture Discriminant Analysis (MDA). It extends LDA and QDA.
- By varying the number of components in each mixture, we can handle classes of any shape, and obtain arbitrarily complex nonlinear decision boundaries.
- We may impose $\Sigma_{kr} = \Sigma$, $\Sigma_{kr} = \sigma_{kr}^2 \mathbf{I}$, or any other parsimonious model, to control the complexity of the model.



Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts



Observed-data likelihood

- Observed-data likelihood:

$$L(\theta) = \prod_{i=1}^n p(x_i, y_i; \theta) = \prod_{i=1}^n p(x_i | y_i; \theta) p(y_i; \theta)$$

$$= \prod_{i=1}^n \prod_{k=1}^c \left(\sum_{r=1}^{R_k} \pi_{kr} \phi(x_i; \mu_{kr}, \Sigma_{kr}) \right)^{y_{ik}} \pi_k^{y_{ik}}$$

- Observed-data log-likelihood:

$$\ell(\theta) = \sum_{k=1}^c \sum_{i=1}^n y_{ik} \log \left(\sum_{r=1}^{R_k} \pi_{kr} \phi(x_i; \mu_{kr}, \Sigma_{kr}) \right) + \sum_{k=1}^c \sum_{i=1}^n y_{ik} \log \pi_k$$

- Again, the EM algorithm can be used to estimate the model parameters (see ESL pp. 399-402 for details).



MDA using package mclust: Iris data

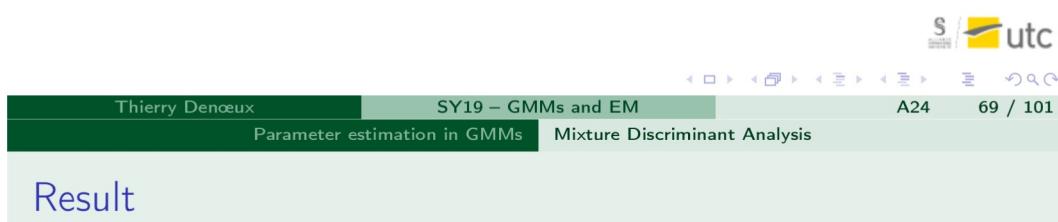
```

odd <- seq(from = 1, to = nrow(iris), by = 2)
even <- odd + 1
X.train <- iris[odd,-5]
Class.train <- iris[odd,5]
X.test <- iris[even,-5]
Class.test <- iris[even,5]

# general covariance structure selected by BIC
irisMclustDA <- MclustDA(X.train, Class.train)
summary(irisMclustDA, newdata = X.test, newclass = Class.test)

plot(irisMclustDA)

```



Result

```
> summary(irisMclustDA, newdata = X.test, newclass = Class.test)
```

Gaussian finite mixture model for classification

MclustDA model summary:

log.likelihood	n	df	BIC
-63.55015	75	53	-355.9272

Classes	n	Model G
setosa	25	VEI 2
versicolor	25	EEV 2
virginica	25	XXX 1

Training classification summary:

Predicted			
Class	setosa	versicolor	virginica
setosa	25	0	0
versicolor	0	25	0
virginica	0	0	25

Training error = 0

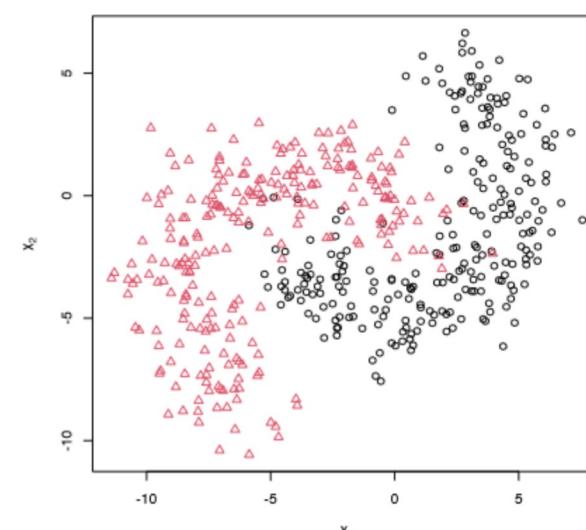
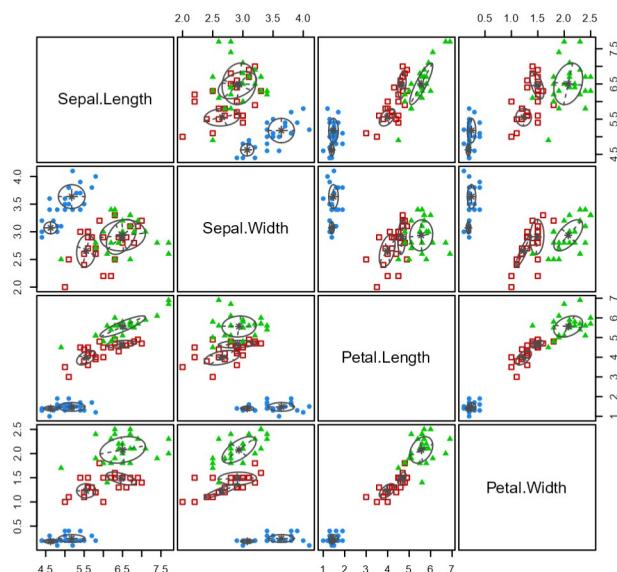
Test classification summary:

Predicted			
Class	setosa	versicolor	virginica
setosa	25	0	0
versicolor	0	24	1
virginica	0	0	25

Test error = 0.01333333



Result



MDA using package mclust: Bananas data

Result

```
> summary(res, newdata = data.test$x, newclass = data.test$y)
-----
Gaussian finite mixture model for classification
-----

McclusDA model summary:

log-likelihood n df      BIC
-2633.035 500 26 -5427.649

Classes n % Model G
 1 250 50   EEV 3
 2 250 50   EEV 3

Training confusion matrix:
  Predicted
Class 1 2
  1 241 9
  2 10 240
Classification error = 0.038
Brier score        = 0.0306

Test confusion matrix:
  Predicted
Class 1 2
  1 471 29
  2 18 482
Classification error = 0.047
Brier score        = 0.0378
```

Thierry Denœux

SY19 – GMMs and EM

Regression models



Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

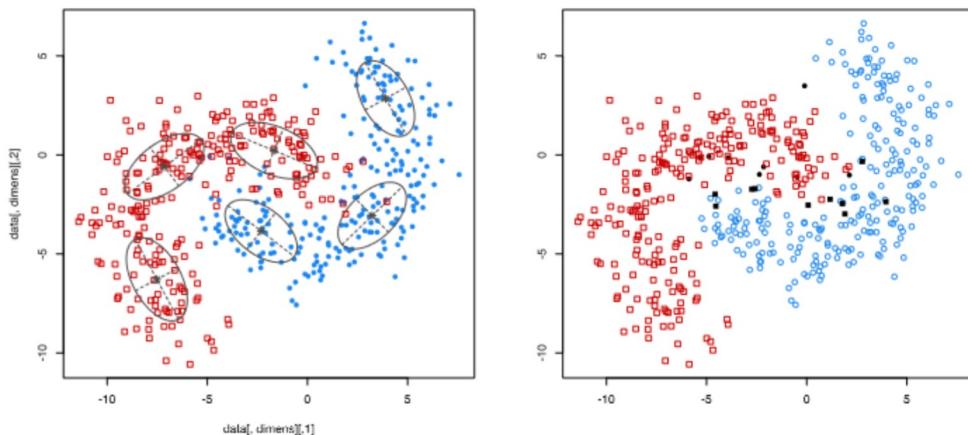
2 Parameter estimation in GMMs

- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts

Result



Thierry Denœux

SY19 – GMMs and EM

Mixture of regressions



Overview

1 EM algorithm

- General formulation
- Simple example
- Analysis

2 Parameter estimation in GMMs

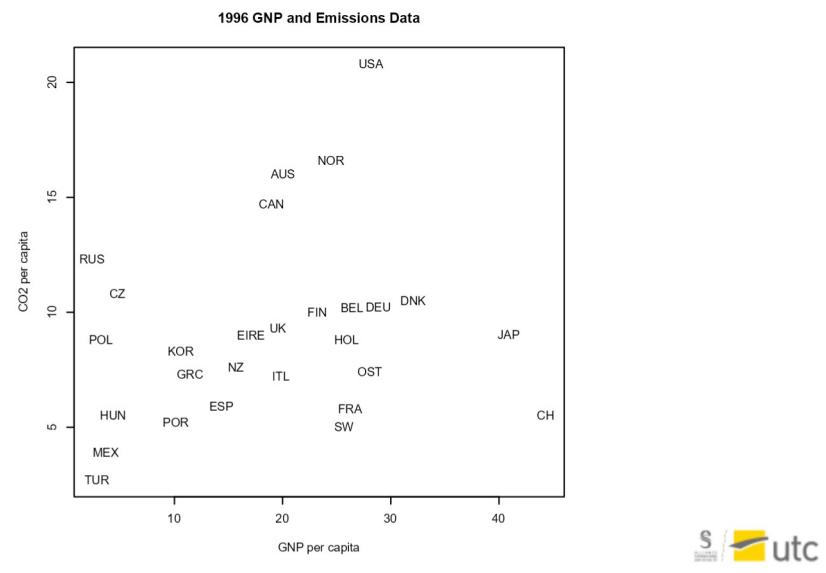
- Unsupervised learning
- Semi-supervised learning
- Mixture Discriminant Analysis

3 Regression models

- Mixture of regressions
- Mixture of experts



Introductory example



Introductory example (continued)

- The data in the previous slide do not show any clear linear trend.
- However, there seem to be several groups for which a linear model would be a reasonable approximation.
- How to identify those groups and the corresponding linear models?

Formalization

- We assume that the response variable Y depends on the input variable X in different ways, depending on a latent variable Z . (Beware: we have switched back to the classical notation for regression models!)
- This model is called **mixture of regressions** or **switching regressions**. It has been widely studied in the econometrics literature.

Model

- Model:

$$Y = \begin{cases} \beta_1^T X + \epsilon_1, & \epsilon_1 \sim \mathcal{N}(0, \sigma_1^2) \quad \text{if } Z = 1 \\ \vdots \\ \beta_c^T X + \epsilon_c, & \epsilon_c \sim \mathcal{N}(0, \sigma_c^2) \quad \text{if } Z = c \end{cases}$$

with $X = (1, X_1, \dots, X_p)$, and

$$\mathbb{P}(Z = k) = \pi_k, \quad k = 1, \dots, c$$

- So, the marginal pdf of Y is

$$p(y | X = x) = \sum_{k=1}^c \pi_k \phi(y; \beta_k^T x, \sigma_k^2)$$

Observed and complete-data likelihoods

- Observed-data likelihood:

$$L(\theta) = \prod_{i=1}^n p(y_i; \theta) = \prod_{i=1}^n \sum_{k=1}^c \pi_k \phi(y_i; \beta_k^T x_i, \sigma_k^2)$$

- Complete-data likelihood:

$$\begin{aligned} L_c(\theta) &= \prod_{i=1}^n p(y_i, z_i; \theta) = \prod_{i=1}^n p(y_i | z_i; \theta) p(z_i | \pi) \\ &= \prod_{i=1}^n \prod_{k=1}^c \phi(y_i; \beta_k^T x_i, \sigma_k^2)^{z_{ik}} \pi_k^{z_{ik}} \end{aligned}$$

with $z_{ik} = I(z_i = k)$.



EM algorithm

E-step: Compute

$$\begin{aligned} z_{ik}^{(t)} &= \mathbb{P}_{\theta^{(t)}}[Z_i = k | y_i] \\ &= \frac{\phi(y_i; \beta_k^{(t)T} x_i, (\sigma_k^{(t)})^2) \pi_k^{(t)}}{\sum_{\ell=1}^c \phi(y_i; \beta_\ell^{(t)T} x_i, (\sigma_\ell^{(t)})^2) \pi_\ell^{(t)}} \end{aligned}$$

M-step: Maximize $Q(\theta, \theta^{(t)})$. As before, we get

$$\pi_k^{(t+1)} = \frac{n_k^{(t)}}{n}$$

with $n_k^{(t)} = \sum_{i=1}^n z_{ik}^{(t)}$.



Derivation of function Q

- Complete-data log-likelihood:

$$\ell_c(\theta) = \sum_{i=1}^n \sum_{k=1}^c z_{ik} \log \phi(y_i; \beta_k^T x_i, \sigma_k^2) + \sum_{i=1}^n \sum_{k=1}^c z_{ik} \log \pi_k$$

- It is linear in the z_{ik} . Consequently, the Q function is simply

$$Q(\theta, \theta^{(t)}) = \underbrace{\sum_{k=1}^c \sum_{i=1}^n z_{ik}^{(t)} \log \phi(y_i; \beta_k^T x_i, \sigma_k^2)}_{\text{term depending on } \beta_k \text{ and } \sigma_k} + \underbrace{\sum_{i=1}^n \sum_{k=1}^c z_{ik}^{(t)} \log \pi_k}_{\text{term depending on } \{\pi_k\}}$$

with $z_{ik}^{(t)} = \mathbb{E}_{\theta^{(t)}}[Z_{ik} | y_i] = \mathbb{P}_{\theta^{(t)}}[Z_i = k | y_i]$.



M-step: update of the β_k and σ_k |

- In $Q(\theta, \theta^{(t)})$, the term depending on β_k is

$$\begin{aligned} \sum_{i=1}^n z_{ik}^{(t)} \log \phi(y_i; \beta_k^T x_i, \sigma_k^2) &= \sum_{i=1}^n z_{ik}^{(t)} \left[-\frac{\log(2\pi\sigma_k^2)}{2} - \frac{1}{2\sigma_k^2} (y_i - \beta_k^T x_i)^2 \right] \\ &= -\frac{1}{2\sigma_k^2} \underbrace{\sum_{i=1}^n z_{ik}^{(t)} (y_i - \beta_k^T x_i)^2}_{SS_k} \\ &\quad - \frac{n_k^{(t)} \log(2\pi\sigma_k^2)}{2} \end{aligned}$$



M-step: update of the β_k and σ_k II

- Minimizing SS_k w.r.t. β_k is a weighted least-squares (WLS) problem.
In matrix form,

$$SS_k = (\mathbf{y} - \mathbf{X}\beta_k)^T \mathbf{W}_k^{(t)} (\mathbf{y} - \mathbf{X}\beta_k)$$

where $\mathbf{W}_k^{(t)} = \text{diag}(z_{1k}^{(t)}, \dots, z_{nk}^{(t)})$ is a diagonal matrix of size n .

- The solution is the WLS estimate of β_k :

$$\beta_k^{(t+1)} = (\mathbf{X}^T \mathbf{W}_k^{(t)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_k^{(t)} \mathbf{y}$$

M-step: update of the β_k and σ_k III

- Plugging in the estimate $\beta_k^{(t+1)}$ in the expression of the Q function and differentiating with respect to σ_k^2 , we obtain the value of σ_k^2 maximizing $Q(\theta, \theta^{(t)})$ as the average of the squared residuals weighted by the $z_{ik}^{(t)}$:

$$\begin{aligned} \sigma_k^{2(t+1)} &= \frac{1}{n_k^{(t)}} \sum_{i=1}^n z_{ik}^{(t)} (y_i - \beta_k^{(t+1)T} x_i)^2 \\ &= \frac{1}{n_k^{(t)}} (\mathbf{y} - \mathbf{X}\beta_k^{(t+1)})^T \mathbf{W}_k^{(t)} (\mathbf{y} - \mathbf{X}\beta_k^{(t+1)}) \end{aligned}$$

Thierry Denœux SY19 – GMMs and EM A24 85 / 101

Regression models Mixture of regressions

Mixture of regressions using mixtools

```
library(mixtools)
data(CO2data)
attach(CO2data)

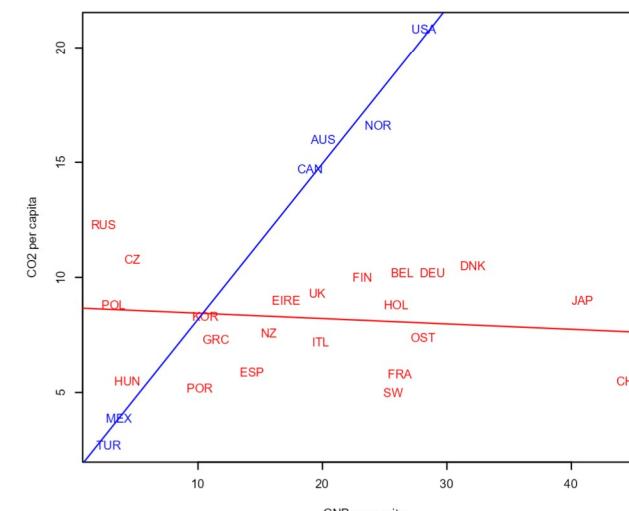
CO2reg <- regmixEM(CO2, GNP)
summary(CO2reg)

ii1<-CO2reg$posterior>0.5
ii2<-CO2reg$posterior<=0.5
text(GNP[ii1],CO2[ii1],country[ii1],col='red')
text(GNP[ii2],CO2[ii2],country[ii2],col='blue')
abline(CO2reg$beta[,1],col='red')
abline(CO2reg$beta[,2],col='blue')
```

Thierry Denœux SY19 – GMMs and EM A24 86 / 101

Regression models Mixture of regressions

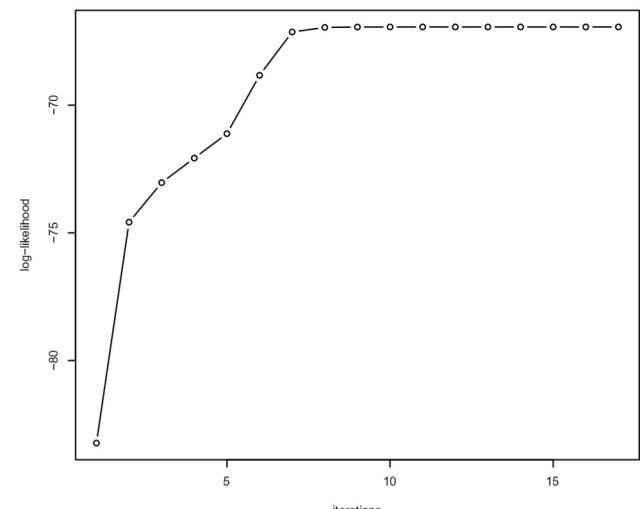
Best solution in 10 runs



Thierry Denœux SY19 – GMMs and EM A24 87 / 101

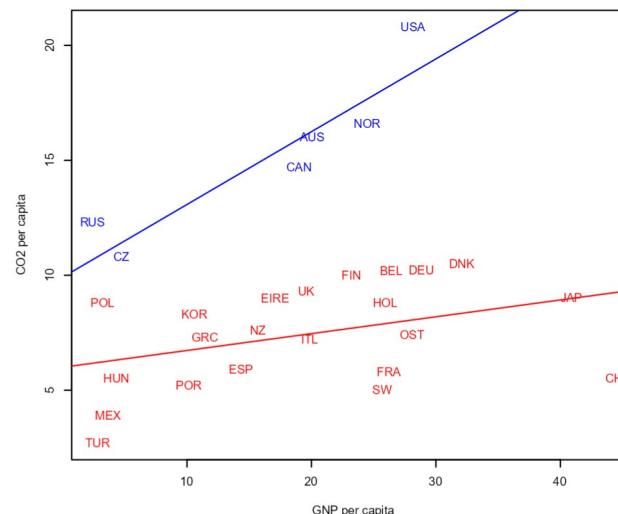
Thierry Denœux SY19 – GMMs and EM A24 88 / 101

Increase of log-likelihood



Increase of log-likelihood

Another solution (with lower log-likelihood)

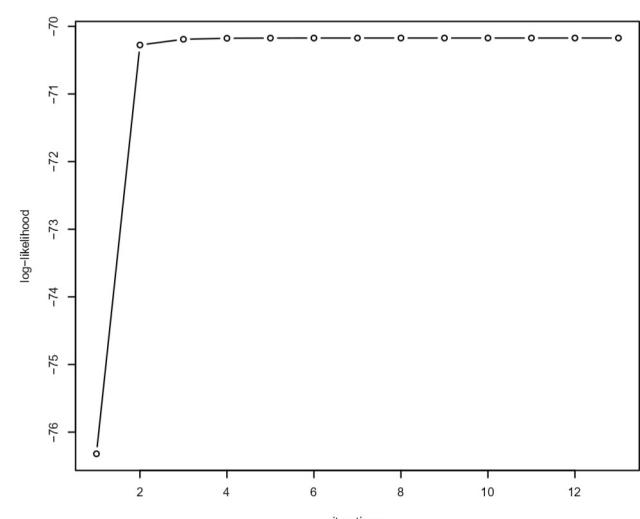


Overview

- ① EM algorithm
 - General formulation
 - Simple example
 - Analysis

 - ② Parameter estimation in GMMs
 - Unsupervised learning
 - Semi-supervised learning
 - Mixture Discriminant Analysis

 - ③ Regression models
 - Mixture of regressions
 - Mixture of experts



Making the mixing proportions predictor-dependent

- An interesting extension of the previous model is to assume the proportions π_k to be partially explained by a vector of **concomitant variables** W .
- If $W = X$, we can approximate the regression function by different linear functions in different regions of the predictor space.
- In ML, this method is referred to as the **mixture of experts** method.
- A useful parametric form for π_k that ensures $\pi_k \geq 0$ and $\sum_{k=1}^c \pi_k = 1$ is the **multinomial logit (softmax)** model:

$$\pi_k(w, \alpha) = \frac{\exp(\alpha_k^T w)}{\sum_{l=1}^c \exp(\alpha_l^T w)}$$

with $\alpha = (\alpha_1, \dots, \alpha_c)$ and $\alpha_1 = 0$.



Generalized EM algorithm

- To ensure the convergence of EM, we only need, at the M step of each iteration t , to find an estimate $\theta^{(t+1)}$ such that

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

- Any algorithm that chooses $\theta^{(t+1)}$ at each iteration to guarantee the above condition without maximizing $Q(\theta, \theta^{(t)})$ is called a **Generalized EM (GEM) algorithm**.
- Here, we can perform a single step of the Newton-Raphson algorithm to maximize

$$\sum_{i=1}^n \sum_{k=1}^c z_{ik}^{(t)} \log \pi_k(w_i, \alpha)$$

with respect to α .

- Backtracking can be used to ensure ascent.



EM algorithm

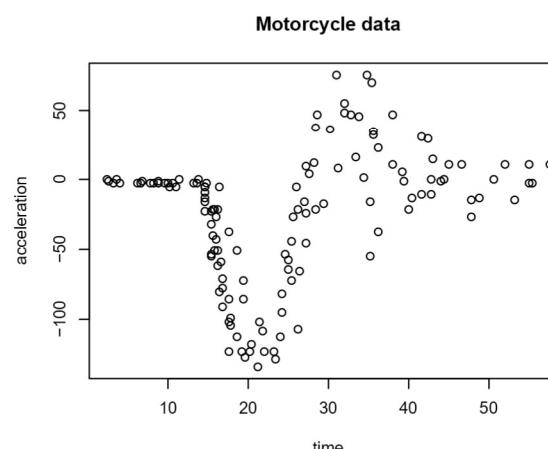
- The Q function is the same as before, except that the π_k now depend on the w_i and parameter α :

$$Q(\theta, \theta^{(t)}) = \sum_{i=1}^n \sum_{k=1}^c z_{ik}^{(t)} \log \phi(y_i; \beta_k^T x_i, \sigma_k^2) + \sum_{i=1}^n \sum_{k=1}^c z_{ik}^{(t)} \log \pi_k(w_i, \alpha)$$

- In the M-step, the update formula for β_k and σ_k are unchanged.
- The last term of $Q(\theta, \theta^{(t)})$ can be maximized w.r.t. α using an iterative algorithm, such as the Newton-Raphson procedure. (See remark on next slide)



Example: motorcycle data

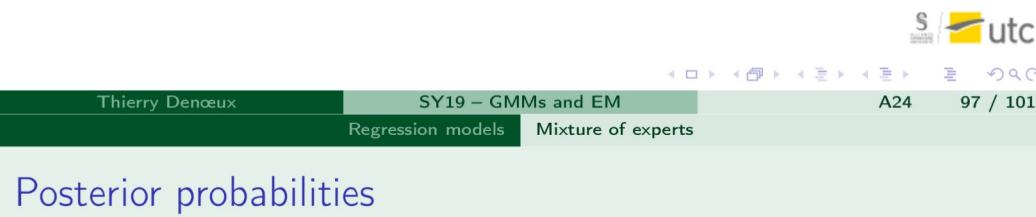


```
library('MASS')
x<-mcycle$times
y<-mcycle$accel
plot(x,y)
```

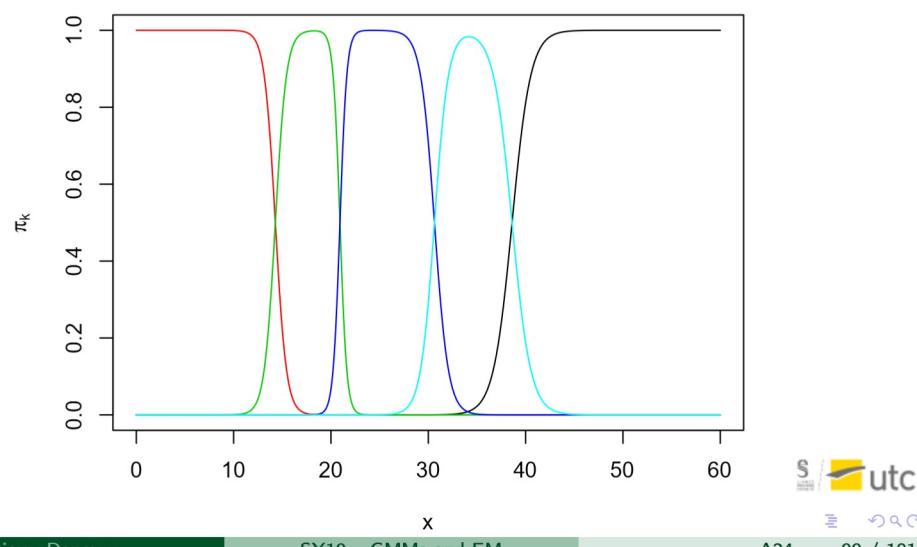


Mixture of experts using flexmix

```
library(flexmix)
K<-5
res<-flexmix(y ~ x,k=K,model=FLXMRglm(family="gaussian"),
concomitant=FLXPmultinom(formula=~x))
beta<- parameters(res)[1:2]
alpha<-res@concomitant@coef
```



Motorcycle data – posterior probabilities



Plotting the posterior probabilities

```
xt<-seq(0,60,0.1)
Nt<-length(xt)
plot(x,y)
pit=matrix(0,Nt,K)
for(k in 1:K) pit[,k]<-exp(alpha[1,k]+alpha[2,k]*xt)
pit<-pit/rowSums(pit)

plot(xt,pit[,1],type="l",col=1)
for(k in 2:K) lines(xt,pit[,k],col=k)
```



```
yhat<-rep(0,Nt)
for(k in 1:K) yhat<-yhat+pit[,k]*(beta[1,k]+beta[2,k]*xt)

plot(x,y,main="Motorcycle data",xlab="time",ylab="acceleration")
for(k in 1:K) abline(beta[1:2,k],lty=2)
lines(xt,yhat,col='red',lwd=2)
```



Regression lines and predictions

