

Computational statistics

EM algorithm

Thierry Denœux

February-March 2017

EM Algorithm

- An iterative optimization strategy motivated by a notion of **missingness** and by consideration of the conditional distribution of what is missing given what is observed.
- Can be very **simple to implement**. Can reliably find an optimum through stable, uphill steps.
- Difficult likelihoods often arise when data are missing. EM simplifies such problems. In fact, the 'missing data' may not truly be missing: they may be only a conceptual ploy to exploit the EM simplification!

Overview

EM algorithm

- Description

- Analysis

Some variants

- Facilitating the E-step

- Facilitating the M-step

Variance estimation

- Louis' method

- SEM algorithm

- Bootstrap

Application to Regression models

- Mixture of regressions

- Mixture of experts

Overview

EM algorithm

- Description

- Analysis

Some variants

- Facilitating the E-step

- Facilitating the M-step

Variance estimation

- Louis' method

- SEM algorithm

- Bootstrap

Application to Regression models

- Mixture of regressions

- Mixture of experts

Notation

\mathbf{Y} : Observed variables.

\mathbf{Z} : Missing or latent variables.

\mathbf{X} : Complete data $\mathbf{X} = (\mathbf{Y}, \mathbf{Z})$.

θ : Unknown parameter.

$L(\theta)$: observed-data likelihood, short for $L(\theta; \mathbf{y}) = f(\mathbf{y}; \theta)$

$L_c(\theta)$: complete-data likelihood, short for $L(\theta; \mathbf{x}) = f(\mathbf{x}; \theta)$

$\ell(\theta), \ell_c(\theta)$: observed and complete-data log-likelihoods.

Notation

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

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) &= \mathbb{E}_{\boldsymbol{\theta}^{(t)}} \{ \ell_c(\boldsymbol{\theta}) \mid \mathbf{y} \} \\ &= \mathbb{E}_{\boldsymbol{\theta}^{(t)}} \{ \log f(\mathbf{X}; \boldsymbol{\theta}) \mid \mathbf{y} \} \\ &= \int [\log f(\mathbf{x}; \boldsymbol{\theta})] f(\mathbf{z} \mid \mathbf{y}; \boldsymbol{\theta}^{(t)}) d\mathbf{z} \end{aligned}$$

where the last equation emphasizes that \mathbf{Z} is the only random part of \mathbf{X} once we are given $\mathbf{Y} = \mathbf{y}$.

The EM Algorithm

Start with $\theta^{(0)}$. 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 unless a stopping criterion has been met; e.g.,

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

Convergence of the EM Algorithm

- It can be proved that $L(\boldsymbol{\theta})$ increases after each EM iteration, i.e., $L(\boldsymbol{\theta}^{(t+1)}) \geq L(\boldsymbol{\theta}^{(t)})$ for $t = 0, 1, \dots$
- Consequently, the algorithm converges to a local maximum of $L(\boldsymbol{\theta})$ if the likelihood function is bounded above.

Mixture of normal and uniform distributions

- Let $\mathbf{Y} = (Y_1, \dots, Y_n)$ be an i.i.d. sample from a mixture of a normal distribution $\mathcal{N}(\mu, \sigma)$ and a uniform distribution $\mathcal{U}([-a, a])$, with pdf

$$f(y; \theta) = \pi\phi(y; \mu, \sigma) + (1 - \pi)c, \quad (1)$$

where $\phi(\cdot; \mu, \sigma)$ is the normal pdf, $c = (2a)^{-1}$, π is the proportion of the normal distribution in the mixture and $\boldsymbol{\theta} = (\mu, \sigma, \pi)^T$ is the vector of parameters.

- Typically, the uniform distribution corresponds to outliers in the data. The proportion of outliers in the population is then $1 - \pi$.
- We want to estimate $\boldsymbol{\theta}$.

Observed and complete-data likelihoods

- Let $Z_i = 1$ if observation i is not an outlier, $Z_i = 0$ otherwise. We have $Z_i \sim \mathcal{B}(\pi)$.
- The vector $\mathbf{Z} = (Z_1, \dots, Z_n)$ is the missing data.
- Observed-data likelihood:

$$L(\boldsymbol{\theta}) = \prod_{i=1}^n f(y_i; \boldsymbol{\theta}) = \prod_{i=1}^n [\pi \phi(y_i; \mu, \sigma) + (1 - \pi)c]$$

- Complete-data likelihood:

$$\begin{aligned} L_c(\boldsymbol{\theta}) &= \prod_{i=1}^n f(y_i, z_i; \boldsymbol{\theta}) = \prod_{i=1}^n f(y_i | z_i; \mu, \sigma) f(z_i | \pi) \\ &= \prod_{i=1}^n [\phi(y_i; \mu, \sigma)^{z_i} c^{1-z_i} \pi^{z_i} (1 - \pi)^{1-z_i}] \end{aligned}$$

Derivation of function Q

- Complete-data log-likelihood:

$$\ell_c(\boldsymbol{\theta}) = \sum_{i=1}^n z_i \log \phi(y_i; \mu, \sigma) + \left(n - \sum_{i=1}^n z_i \right) \log c + \sum_{i=1}^n (z_i \log \pi + (1 - z_i) \log(1 - \pi))$$

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

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \sum_{i=1}^n z_i^{(t)} \log \phi(y_i; \mu, \sigma) + \left(n - \sum_{i=1}^n z_i^{(t)} \right) \log c + \sum_{i=1}^n \left(z_i^{(t)} \log \pi + (1 - z_i^{(t)}) \log(1 - \pi) \right)$$

with $z_i^{(t)} = \mathbb{E}_{\boldsymbol{\theta}^{(t)}}[Z_i | y_i]$.

EM algorithm

E-step: compute

$$\begin{aligned} z_i^{(t)} &= \mathbb{E}_{\theta^{(t)}}[Z_i | y_i] = \mathbb{P}_{\theta^{(t)}}[Z_i = 1 | y_i] \\ &= \frac{\phi(y_i; \mu^{(t)}, \sigma^{(t)})\pi^{(t)}}{\phi(y_i; \mu^{(t)}, \sigma^{(t)})\pi^{(t)} + c(1 - \pi^{(t)})} \end{aligned}$$

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

$$\pi^{(t+1)} = \frac{1}{n} \sum_{i=1}^n z_i^{(t)}, \quad \mu^{(t+1)} = \frac{\sum_{i=1}^n z_i^{(t)} y_i}{\sum_{i=1}^n z_i^{(t)}}$$

$$\sigma^{(t+1)} = \sqrt{\frac{\sum_{i=1}^n z_i^{(t)} (y_i - \mu^{(t+1)})^2}{\sum_{i=1}^n z_i^{(t)}}}$$

Remark

- As mentioned before, the EM algorithm finds only a local maximum of $\ell(\theta)$.
- It is easy to find a global maximum: if μ is equal to some y_i and $\sigma = 0$, then $\phi(y_i; \mu, \sigma) = \infty$ and, consequently, $\ell(\theta) = +\infty$.
- We are not interested in these global maxima, because they correspond to **degenerate solutions!**

Bayesian posterior mode

- Consider a Bayesian estimation problem with likelihood $L(\theta)$ and prior $f(\theta)$.
- The posterior density is proportional to $L(\theta)f(\theta)$. It can also be maximized by the EM algorithm.
- The E-step requires

$$Q(\theta, \theta^{(t)}) = \mathbb{E}_{\theta^{(t)}} \{ \ell_c(\theta) \mid \mathbf{y} \} + \log f(\theta)$$

- The addition of the log-prior often makes it more difficult to maximize Q during the M-step.
- Some methods can be used to facilitate the M-step in difficult situations (see below).

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to Regression models

Mixture of regressions

Mixture of experts

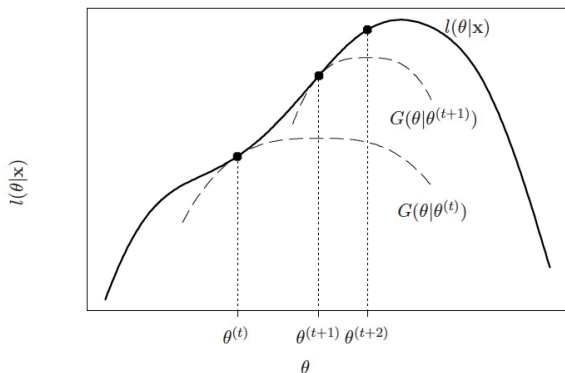
Why does it work?

- **Ascent:** Each M-step increases the log likelihood.
- **Optimization transfer:**

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

- The last two terms in $G(\theta, \theta^{(t)})$ are constant with respect to θ , so Q and G are maximized at the same θ .
- Further, G is tangent to ℓ at $\theta^{(t)}$, and lies everywhere below ℓ . We say that G is a **minorizing function** for ℓ .
- EM transfers optimization from ℓ to the surrogate function G , which is more convenient to maximize.

The nature of EM



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

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

- Consequently,

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

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

$$\ell(\theta) = Q(\theta, \theta^{(t)}) - \underbrace{\mathbb{E}_{\theta^{(t)}}[\log f(Z|y; \theta)|y]}_{H(\theta, \theta^{(t)})} \quad (2)$$

Proof - the minorizing function

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

$$H(\theta, \theta^{(t)}) - H(\theta^{(t)}, \theta^{(t)}) = \mathbb{E}_{\theta^{(t)}} \left[\log \frac{f(Z|y; \theta)}{f(Z|y; \theta^{(t)})} | y \right] \quad (3a)$$

$$\leq \log \mathbb{E}_{\theta^{(t)}} \left[\frac{f(Z|y; \theta)}{f(Z|y; \theta^{(t)})} | y \right] (*) \quad (3b)$$

$$= \log \int f(z|y; \theta) dz = 0 \quad (3c)$$

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

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

$$H(\theta, \theta^{(t)}) \leq H(\theta^{(t)}, \theta^{(t)}) = Q(\theta^{(t)}, \theta^{(t)}) - \ell(\theta^{(t)}), \text{ or}$$

$$\ell(\theta) \geq Q(\theta, \theta^{(t)}) + \ell(\theta^{(t)}) - Q(\theta^{(t)}, \theta^{(t)}) = G(\theta, \theta^{(t)}) \quad (4)$$

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

- From (4), $\ell(\theta^{(t)}) = G(\theta^{(t)}, \theta^{(t)})$.
- Now, we can rewrite (4) as

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

Consequently, $\theta^{(t)}$ maximizes $Q(\theta, \theta^{(t)}) - \ell(\theta)$, hence

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

and

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

Proof - monotonicity

- From (2),

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

- $A \geq 0$ because $\theta^{(t+1)}$ is a maximizer of $Q(\theta, \theta^{(t)})$, and $B \leq 0$ because, from (3), $\theta^{(t)}$ is a maximizer of $H(\theta, \theta^{(t)})$.
- Hence,

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

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to Regression models

Mixture of regressions

Mixture of experts

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to Regression models

Mixture of regressions

Mixture of experts

Monte Carlo EM (MCEM)

- Replace the t th E step with
 - 1 Draw missing datasets $\mathbf{Z}_1^{(t)}, \dots, \mathbf{Z}_{m^{(t)}}^{(t)}$ i.i.d. from $f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta}^{(t)})$. Each $\mathbf{Z}_j^{(t)}$ is a vector of all the missing values needed to complete the observed dataset, so $\mathbf{X}_j^{(t)} = (\mathbf{y}, \mathbf{Z}_j^{(t)})$ denotes a completed dataset where the missing values have been replaced by $\mathbf{Z}_j^{(t)}$.
 - 2 Calculate $\hat{Q}^{(t+1)}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \frac{1}{m^{(t)}} \sum_{j=1}^{m^{(t)}} \log f(\mathbf{X}_j^{(t)}; \boldsymbol{\theta})$.
- Then $\hat{Q}^{(t+1)}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$ is a Monte Carlo estimate of $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$.
- The M step is modified to maximize $\hat{Q}^{(t+1)}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$.
- Increase $m^{(t)}$ as iterations progress to reduce the Monte Carlo variability of \hat{Q} . MCEM will not converge in the same sense as ordinary EM, rather values of $\boldsymbol{\theta}^{(t)}$ will bounce around the true maximum, with a precision that depends on $m^{(t)}$.

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to Regression models

Mixture of regressions

Mixture of experts

Generalized EM (GEM) algorithm

- In the original EM algorithm, $\theta^{(t+1)}$ is a maximizer of $Q(\theta, \theta^{(t)})$, i.e.,

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

for all θ .

- However, to ensure convergence, we only need 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**.

EM gradient algorithm

- Replace the M step with a single step of Newton's method, thereby approximating the maximum without actually solving for it exactly.
- Instead of maximizing, choose:

$$\begin{aligned}\boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} - \mathbf{Q}''(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})^{-1} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \mathbf{Q}'(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \\ &= \boldsymbol{\theta}^{(t)} - \mathbf{Q}''(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})^{-1} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \ell'(\boldsymbol{\theta}^{(t)})\end{aligned}$$

- Ascent is ensured for canonical parameters in exponential families. Backtracking can ensure ascent in other cases; inflating steps can speed convergence.

ECM algorithm

- Replaces the M step with a series of computationally simpler conditional maximization (CM) steps.
- Call the collection of simpler CM steps after the t th E step a CM **cycle**. Thus, the t th iteration of ECM is comprised of the t th E step and the t th CM cycle.
- Let S denote the total number of CM steps in each CM cycle.

ECM algorithm (continued)

- For $s = 1, \dots, S$, the s th CM step in the t th cycle requires the maximization of $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$ subject to (or conditional on) a constraint, say

$$\mathbf{g}_s(\boldsymbol{\theta}) = \mathbf{g}_s(\boldsymbol{\theta}^{(t+(s-1)/S)})$$

where $\boldsymbol{\theta}^{(t+(s-1)/S)}$ is the maximizer found in the $(s-1)$ th CM step of the current cycle.

- When the entire cycle of S steps of CM has been completed, we set $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t+S/S)}$ and proceed to the E step for the $(t+1)$ th iteration.
- ECM is a GEM algorithm, since each CM step increases Q .
- The art of constructing an effective ECM algorithm lies in choosing the constraints cleverly.

Choice 1: Iterated Conditional Modes / Gauss-Seidel

- Partition θ into S subvectors, $\theta = (\theta_1, \dots, \theta_S)$.
- In the s th CM step, maximize Q with respect to θ_s while holding all other components of θ fixed.
- This amounts to the constraint induced by the function

$$g_s(\theta) = (\theta_1, \dots, \theta_{s-1}, \theta_{s+1}, \dots, \theta_S).$$

Choice 2

- At the s th CM step, maximize Q with respect to all other components of θ while holding θ_s fixed.
- Then $g_s(\theta) = \theta_s$.
- Additional systems of constraints can be imagined, depending on the particular problem context.
- A variant of ECM inserts an E step between each pair of CM steps, thereby updating Q at every stage of the CM cycle.

Overview

EM algorithm

- Description

- Analysis

Some variants

- Facilitating the E-step

- Facilitating the M-step

Variance estimation

- Louis' method

- SEM algorithm

- Bootstrap

Application to Regression models

- Mixture of regressions

- Mixture of experts

Variance of the MLE

- Let $\hat{\theta}$ be the MLE of θ .
- As $n \rightarrow \infty$, the limiting distribution of $\hat{\theta}$ is $\mathcal{N}(\theta^*, I(\theta^*)^{-1})$, where θ^* is the true value of θ , and

$$I(\theta) = \mathbb{E}[\ell'(\theta)\ell'(\theta)^T] = -\mathbb{E}[\ell''(\theta)]$$

is the **expected Fisher information matrix** (the second equality holds under some regularity conditions).

- $I(\theta^*)$ can be estimated by $I(\hat{\theta})$, or by $-\ell''(\hat{\theta}) = I_{obs}(\hat{\theta})$ (**observed information matrix**).
- Standard error estimates can be obtained by computing the square roots of the diagonal elements of $I_{obs}(\hat{\theta})^{-1}$.

Obtaining variance estimates

- The EM algorithm allows us to estimate $\hat{\theta}$, but it does not directly provide an estimate of $I(\theta^*)$.
- Direct computation of $I(\hat{\theta})$ or $I_{obs}(\hat{\theta})$ is often difficult.
- Main methods:
 - 1 Louis' method
 - 2 Supplement EM (SEM) algorithm
 - 3 Bootstrap

Overview

EM algorithm

- Description

- Analysis

Some variants

- Facilitating the E-step

- Facilitating the M-step

Variance estimation

- Louis' method

- SEM algorithm

- Bootstrap

Application to Regression models

- Mixture of regressions

- Mixture of experts

Missing information principle

- We have seen that

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

from which we get

$$\ell(\boldsymbol{\theta}) = \ell_c(\boldsymbol{\theta}) - \log f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta}).$$

- Differentiating twice and negating both sides, then taking expectations over the conditional distribution of \mathbf{X} given \mathbf{y} ,

$$\underbrace{-\ell''(\boldsymbol{\theta})}_{\hat{\mathbf{i}}_{\mathbf{Y}}(\boldsymbol{\theta})} = \underbrace{\mathbb{E}[-\ell''_c(\boldsymbol{\theta})|\mathbf{y}]}_{\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta})} - \underbrace{\mathbb{E}\left[-\frac{\partial^2 \log f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \middle| \mathbf{y}\right]}_{\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta})}$$

where

- $\hat{\mathbf{i}}_{\mathbf{Y}}(\boldsymbol{\theta})$ is the **observed information**,
- $\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta})$ is the **complete information**, and
- $\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta})$ is the **missing information**.

Louis' method

- Computing $\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta})$ and $\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta})$ is sometimes easier than computing $-\ell''(\boldsymbol{\theta})$ directly
- We can show that

$$\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta}) = \text{Var}[S_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta})],$$

where the variance is taken w.r.t. $\mathbf{Z}|\mathbf{y}$, and

$$S_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta}) = \frac{\partial f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

is the conditional score.

- As the expected score is zero at $\hat{\boldsymbol{\theta}}$, we have

$$\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\hat{\boldsymbol{\theta}}) = \int S_{\mathbf{Z}|\mathbf{Y}}(\hat{\boldsymbol{\theta}}) S_{\mathbf{Z}|\mathbf{Y}}(\hat{\boldsymbol{\theta}})^T f(\mathbf{z}|\mathbf{y}; \hat{\boldsymbol{\theta}}) d\mathbf{z}$$

Monte Carlo approximation

- When they cannot be computed analytically, $\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta})$ and $\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta})$ can sometimes be approximated by Monte Carlo simulation.
- Method: generate simulated datasets $\mathbf{x}_j = (\mathbf{y}, \mathbf{z}_j)$, $j = 1, \dots, N$, where \mathbf{y} is the observed dataset, and the \mathbf{z}_j are imputed missing datasets drawn from $f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta})$
- Then,

$$\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta}) \approx \frac{1}{N} \sum_{j=1}^N -\frac{\partial^2 \log f(\mathbf{x}_j; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}$$

and $\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta})$ is approximated by the sample variance of the values

$$\frac{\partial f(\mathbf{z}_j|\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to Regression models

Mixture of regressions

Mixture of experts

EM mapping

- Let Ψ denotes the EM mapping, defined by

$$\theta^{(t+1)} = \Psi(\theta^{(t)})$$

- From the convergence of EM, $\hat{\theta}$ is a fixed point:

$$\hat{\theta} = \Psi(\hat{\theta}).$$

- The Jacobian matrix of Ψ is the $p \times p$ matrix

$$\Psi'(\theta) = \left(\frac{\partial \Psi_i(\theta)}{\partial \theta_j} \right).$$

- It can be shown that

$$\Psi'(\hat{\theta})^T = \hat{i}_{Z|Y}(\hat{\theta}) \hat{i}_X(\hat{\theta})^{-1}$$

Using $\Psi'(\theta)$ for variance estimation

- From the missing information principle,

$$\begin{aligned}\hat{\mathbf{i}}_{\mathbf{Y}}(\hat{\theta}) &= \hat{\mathbf{i}}_{\mathbf{X}}(\hat{\theta}) - \hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\hat{\theta}) \\ &= \left[\mathbf{I} - \hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\hat{\theta})\hat{\mathbf{i}}_{\mathbf{X}}(\hat{\theta})^{-1} \right] \hat{\mathbf{i}}_{\mathbf{X}}(\hat{\theta}) \\ &= \left[\mathbf{I} - \Psi'(\hat{\theta})^T \right] \hat{\mathbf{i}}_{\mathbf{X}}(\hat{\theta}).\end{aligned}$$

Hence,

$$\hat{\mathbf{i}}_{\mathbf{Y}}(\hat{\theta})^{-1} = \hat{\mathbf{i}}_{\mathbf{X}}(\hat{\theta})^{-1} \left[\mathbf{I} - \Psi'(\hat{\theta})^T \right]^{-1}$$

- From the equality

$$(\mathbf{I} - \mathbf{P})^{-1} = (\mathbf{I} - \mathbf{P} + \mathbf{P})(\mathbf{I} - \mathbf{P})^{-1} = \mathbf{I} + \mathbf{P}(\mathbf{I} - \mathbf{P})^{-1},$$

we get

$$\hat{\mathbf{i}}_{\mathbf{Y}}(\hat{\theta})^{-1} = \hat{\mathbf{i}}_{\mathbf{X}}(\hat{\theta})^{-1} \left\{ \mathbf{I} + \Psi'(\hat{\theta})^T \left[\mathbf{I} - \Psi'(\hat{\theta})^T \right]^{-1} \right\}. \quad (5)$$

Estimation of $\Psi'(\hat{\theta})$

- Let r_{ij} be the element (i, j) of $\Psi'(\hat{\theta})$. By definition,

$$\begin{aligned} r_{ij} &= \frac{\partial \Psi_i(\hat{\theta})}{\partial \theta_j} \\ &= \lim_{\theta_j \rightarrow \hat{\theta}_j} \frac{\Psi_i(\hat{\theta}_1, \dots, \hat{\theta}_{j-1}, \theta_j, \hat{\theta}_{j+1}, \dots, \hat{\theta}_p) - \Psi_i(\hat{\theta})}{\theta_j - \hat{\theta}_j} \\ &= \lim_{t \rightarrow \infty} \frac{\Psi_i(\theta^{(t)}(j)) - \Psi_i(\hat{\theta})}{\theta_j^{(t)} - \hat{\theta}_j} = \lim_{t \rightarrow \infty} r_{ij}^{(t)} \end{aligned}$$

where $\theta^{(t)}(j) = (\hat{\theta}_1, \dots, \hat{\theta}_{j-1}, \theta_j^{(t)}, \hat{\theta}_{j+1}, \dots, \hat{\theta}_p)$, and $(\theta_j^{(t)})$, $t = 1, 2, \dots$ is a sequence of values converging to $\hat{\theta}_j$.

- Method: compute the $r_{ij}^{(t)}$, $t = 1, 2, \dots$ until they stabilize to some values. Then compute $\hat{\mathbf{i}}_{\Psi}(\hat{\theta})^{-1}$ using (5).

SEM algorithm

- 1 Run the EM algorithm to convergence, finding $\hat{\theta}$.
- 2 Restart the algorithm from some $\theta^{(0)}$ near $\hat{\theta}$. For $t = 0, 1, 2, \dots$
 - 1 Take a standard E step and M step to produce $\theta^{(t+1)}$ from $\theta^{(t)}$.
 - 2 For $j = 1, \dots, p$, define $\theta^{(t)}(j) = (\hat{\theta}_1, \dots, \hat{\theta}_{j-1}, \theta_j^{(t)}, \hat{\theta}_{j+1}, \dots, \hat{\theta}_p)$ and

$$r_{ij}^{(t)} = \frac{\psi_i(\theta^{(t)}(j)) - \hat{\theta}_i}{\theta_j^{(t)} - \hat{\theta}_j}$$

for $i = 1, \dots, p$. (Recall that $\Psi(\hat{\theta}) = \hat{\theta}$.)

- 3 Stop when all $r_{ij}^{(t)}$ have converged
- 3 The (i, j) th element of $\Psi'(\hat{\theta})$ equals $\lim_{t \rightarrow \infty} r_{ij}^{(t)}$. Use the final estimate of $\Psi'(\hat{\theta})$ to get the variance.
- 4 SEM is numerically stable and requires little extra work.

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to Regression models

Mixture of regressions

Mixture of experts

Principle

- Consider the case of iid data $\mathbf{y} = (\mathbf{w}_1, \dots, \mathbf{w}_n)$
- If we knew the distribution of the \mathbf{W}_i , we could
 - generate many samples $\mathbf{y}_1, \dots, \mathbf{y}_n$,
 - compute the ML estimate $\hat{\theta}_j$ of θ from each sample \mathbf{y}_j , and
 - estimate the variance of $\hat{\theta}$ by the sample variance of the estimates $\hat{\theta}_1, \dots, \hat{\theta}_N$.
- Bootstrap principle: use the **empirical distribution** in place of the true distribution of the \mathbf{W}_i

Algorithm

- ① Calculate $\hat{\theta}_{EM}$ using a suitable EM approach applied to $\mathbf{y} = (\mathbf{w}_1, \dots, \mathbf{w}_n)$. Let $j = 1$ and set $\hat{\theta}_j^* = \hat{\theta}_{EM}$.
- ② Increment j . Sample pseudo-data $\mathbf{y}_j^* = (\mathbf{w}_{j1}^*, \dots, \mathbf{w}_{jn}^*)$ at random from $(\mathbf{w}_1, \dots, \mathbf{w}_n)$ with replacement.
- ③ Calculate $\hat{\theta}_j^*$ by applying the same EM approach to the pseudo-data \mathbf{y}_j^*
- ④ Stop if $j = B$ (typically, $B \geq 1000$); otherwise return to step 2.

The collection of parameter estimates $\hat{\theta}_1^*, \dots, \hat{\theta}_B^*$ can be used to estimate the variance of $\hat{\theta}$,

$$\widehat{\text{Var}}(\hat{\theta}) = \frac{1}{B} \sum_{j=1}^B (\hat{\theta}_j^* - \overline{\hat{\theta}^*}) (\hat{\theta}_j^* - \overline{\hat{\theta}^*})^T,$$

where $\overline{\hat{\theta}^*}$ is the sample mean of $\hat{\theta}_1^*, \dots, \hat{\theta}_B^*$.

Pros and cons of the bootstrap

1 Advantages:

- The method is very general, complex analytical derivations are avoided.
- Allows the estimation of other aspects of the sampling distribution of $\widehat{\theta}$, such as expectation (bias), quantiles, etc.

- ## 2 Drawback: bootstrap embeds the EM loop in a second loop of B iterations. May be computationally burdensome when the EM algorithm is slow (because, e.g., of a high proportion of missing data, or high dimensionality).

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to Regression models

Mixture of regressions

Mixture of experts

Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

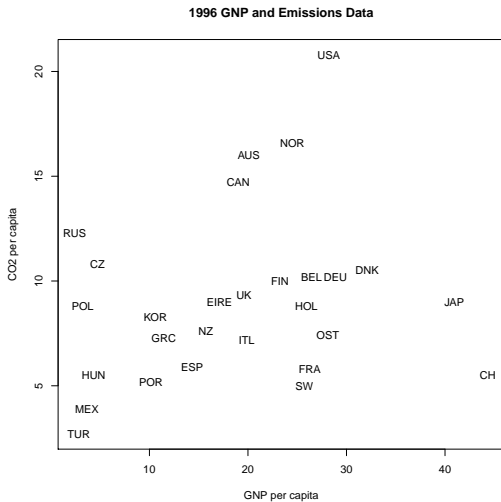
Bootstrap

Application to 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?

Model

- Model: 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:

$$Y = \begin{cases} \beta_1^T X + \epsilon_1, \epsilon_1 \sim \mathcal{N}(0, \sigma_1) & \text{if } Z = 1, \\ \vdots \\ \beta_K^T X + \epsilon_K, \epsilon_K \sim \mathcal{N}(0, \sigma_K) & \text{if } Z = K. \end{cases}$$

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

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

Observed and complete-data likelihoods

- Observed-data likelihood:

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

- Complete-data likelihood:

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

with $z_{ik} = 1$ if $z_i = k$ and $z_{ik} = 0$ otherwise.

Derivation of function Q

- Complete-data log-likelihood:

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

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

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

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

EM algorithm

- E-step: compute

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

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

$$\pi_k^{(t+1)} = \frac{N_k^{(t)}}{N},$$

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

M-step: update of the β_k and σ_k

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

$$SS_k = \sum_{i=1}^N z_{ik}^{(t)} (y_i - \beta_k^T x_i)^2.$$

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

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

with $\mathbf{W}_k = \text{diag}(z_{i1}^{(t)}, \dots, z_{iK}^{(t)})$.

M-step: update of the β_k and σ_k (continued)

- The solution is the WLS estimate of β_k :

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

- The value of σ^k minimizing $Q(\theta, \theta^{(t)})$ is the weighted average of the residuals,

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

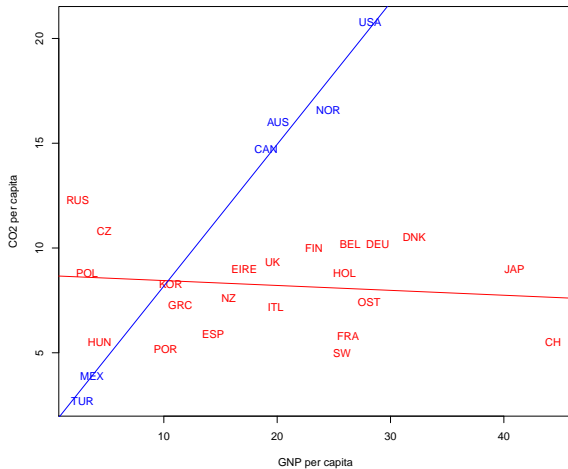
Mixture of regressions using mixtools

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

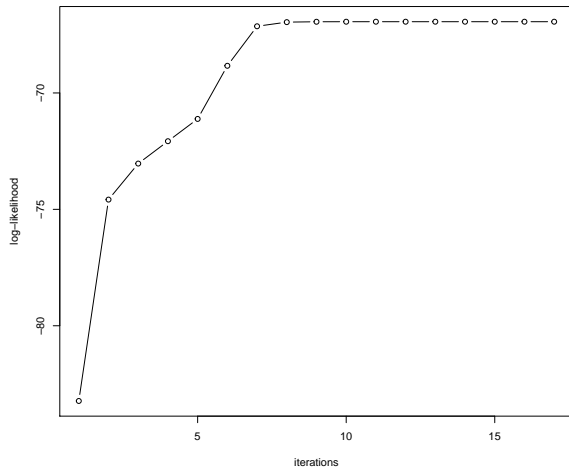
CO2reg <- regmixEM(CO2, GNP)
summary(CO2reg)

ii1<-CO2reg$posterior>0.5
ii2<-CO2reg$posterior<=0.5
text(GNP[ii1],CO2[ii1],country[ii1],col='red')
text(GNP[Cii2],CO2[ii2],country[ii2],col='blue')
abline(CO2reg$beta[,1],col='red')
abline(CO2reg$beta[,2],col='blue')
```

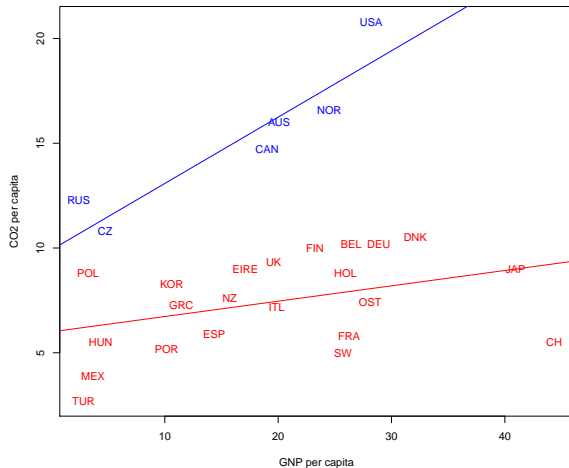
Best solution in 10 runs



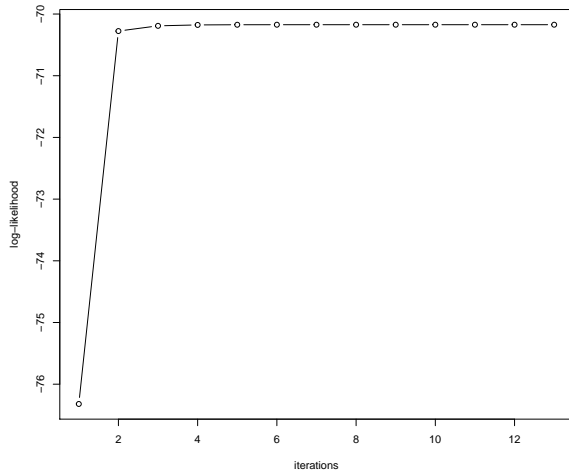
Increase of log-likelihood



Another solution (with lower log-likelihood)



Increase of log-likelihood



Overview

EM algorithm

Description

Analysis

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm

Bootstrap

Application to 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** methods.
- A useful parametric form for π_k that ensures $\pi_k \geq 0$ and $\sum_{k=1}^K \pi_k = 1$ is the multinomial logit model

$$\pi_k(w, \alpha) = \frac{\exp(\alpha_k^T w)}{\sum_{\ell=1}^K \exp(\alpha_\ell^T w)}$$

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

EM algorithm

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

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

- 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)

Generalized EM algorithm

- To ensure convergence of EM, we only need to increase (but not necessarily maximize) $Q(\theta, \theta^{(t)})$ at each step.
- 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

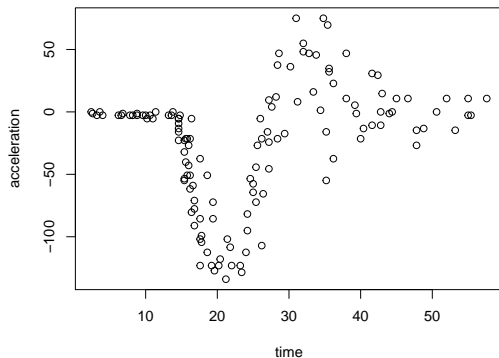
$$\sum_{i=1}^N \sum_{k=1}^K z_{ik}^{(t)} \log \pi_k(w_i, \alpha)$$

with respect to α .

- Backtracking can be used to ensure ascent.

Example: motorcycle data

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
```

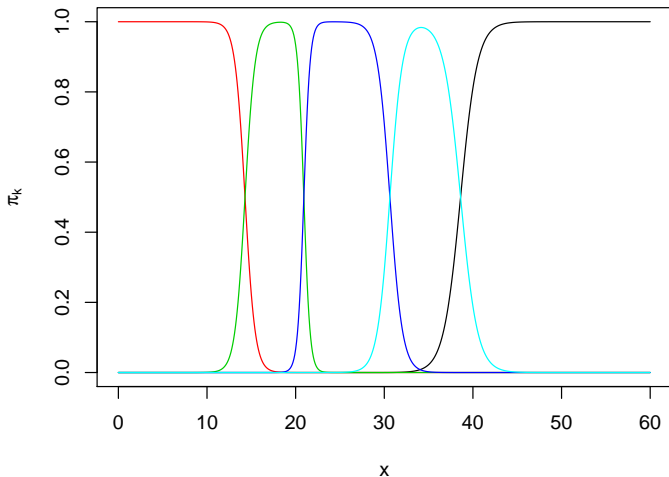
Plotting the posterior probabilities

```
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)
```

Posterior probabilities

Motorcycle data – posterior probabilities



Plotting the predictions

```
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

