

Computational statistics

Chapter 3: EM algorithm

Thierry Denœux
Université de technologie de Compiègne

2022-2023



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 econometrics.
- Can be very simple to implement. Can reliably find an optimum through stable, uphill steps.



Overview

1 EM algorithm

- Description
- Analysis

2 Some variants

- Facilitating the E-step
- Facilitating the M-step

3 Variance estimation

- Louis' method
- SEM algorithm



Overview

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



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



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, conditional on the observed data $\mathbf{Y} = \mathbf{y}$** . Namely

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

$(f(\mathbf{x} \mid \mathbf{y}; \theta^{(t)}) = f(\mathbf{z} \mid \mathbf{y}; \theta^{(t)})$ because \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 Increment t and return to the E step unless a stopping criterion has been met; e.g.,

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

or

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



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.
- Typically, we run the algorithm several times with random initial conditions, and we keep the results of the best run.



Example: 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

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

where $\phi(\cdot; \mu, \sigma)$ is the normal pdf, $c = (2a)^{-1}$ is a known constant, 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:

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



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

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

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



Bayesian posterior mode

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

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

- 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

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



Why does it work?

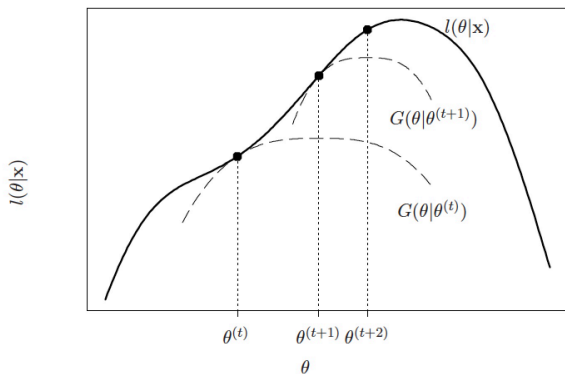
- **Ascent:** Each M-step increases the log likelihood.
- **Optimization transfer:**

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

- The last two terms in $G(\theta, \theta^{(t)})$ do not depend on 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 (see next slide).
- 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 and each M step maximizes it to provide an uphill step.



Proof

- We have

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

- Consequently,

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

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

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



Proof - the minorizing function

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

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

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

$$\leq \log \underbrace{\int f(\mathbf{z} | \mathbf{y}; \theta) d\mathbf{z}}_1 = 0 \quad (3c)$$

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

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



Proof - the minorizing function (continued)

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

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

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

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



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

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



Proof - monotonicity

- From (2),

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

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

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



Overview

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



Overview

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



Monte Carlo EM (MCEM)

- Sometimes, the conditional expectation of $\ell_c(\boldsymbol{\theta})$ given \mathbf{y} cannot be easily computed analytically in the E step.
- Approach: randomly generate sets of missing values according to the conditional distribution $f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta}^{(t)})$, and replace the expectation by an average over generated data sets.



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

$$\widehat{Q}^{(t+1)}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \frac{1}{m^{(t)}} \sum_{j=1}^{m^{(t)}} \log f(\mathbf{X}_j^{(t)}; \boldsymbol{\theta}).$$

- Then $\widehat{Q}^{(t+1)}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$ is a Monte Carlo estimate of $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$.
- The M step is modified to maximize $\widehat{Q}^{(t+1)}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$.



Remarks

- It is advised to increase $m^{(t)}$ as iterations progress to reduce the Monte Carlo variability of \widehat{Q} .
- MCEM will not converge in the same sense as ordinary EM, rather values of $\theta^{(t)}$ will bounce around the true maximum, with a precision that depends on $m^{(t)}$.



Overview

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



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:

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

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



Overview

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



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

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

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 Supplemented EM (SEM) algorithm
 - 3 Bootstrap (to be studied in Chapter 6)



Overview

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



Missing information principle

- We have seen that

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

from which we get

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

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

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

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

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

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

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

is the conditional score.

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

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



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

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



Overview

- 1 EM algorithm
 - Description
 - Analysis
- 2 Some variants
 - Facilitating the E-step
 - Facilitating the M-step
- 3 Variance estimation
 - Louis' method
 - SEM algorithm



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,

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

- Hence,

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



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

- 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

$$\hat{\mathbf{i}}_{\mathbf{Y}}(\hat{\theta})^{-1} = \hat{\mathbf{i}}_{\mathbf{X}}(\hat{\theta})^{-1} \left\{ \mathbf{I} + \Psi'(\hat{\theta})^T \left[\mathbf{I} - \Psi'(\hat{\theta})^T \right]^{-1} \right\} \quad (4)$$

- This result is appealing in that it expresses the desired covariance matrix as the **complete-data covariance matrix** plus an incremental matrix that takes account of **the uncertainty attributable to the missing data**.



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 (4).



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 treating it as the current estimate of θ , run one iteration of EM to obtain $\Psi(\theta^{(t)}(j))$.
 - Obtain the ratio

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

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.

