

# Computational statistics

## Lecture 5: EM algorithm

Thierry Denœux

24 March, 2016



# 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

### Some variants

Facilitating the E-step

Facilitating the M-step

### Variance estimation

Louis' method

SEM algorithm



# Notation

**Y** : Observed variables.

**Z** : Missing or latent variables.

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

$\theta$  : 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

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

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  $\theta$ . 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.



# Trivial example

- $Y, Z$  iid from  $\mathcal{E}(\theta)$  with  $y = 5$  observed but  $z$  missing.
- The complete-data log likelihood function is

$$\ell_c(\theta) = \log\{f_{\mathbf{X}}(\mathbf{x}; \theta)\} = 2 \log(\theta) - \theta y - \theta z.$$

- Thus

$$Q(\theta, \theta^{(t)}) = 2 \log(\theta) - 5\theta - \theta/\theta^{(t)}$$

since  $\mathbb{E}_{\theta^{(t)}}\{Z|y\} = \mathbb{E}_{\theta^{(t)}}\{Z\} = 1/\theta^{(t)}$  follows from independence.

- The maximizer of  $Q(\theta, \theta^{(t)})$  is the root of  $2/\theta - 5 - 1/\theta^{(t)} = 0$ . Thus  $\theta^{(t+1)} = \frac{2\theta^{(t)}}{5\theta^{(t)}+1}$ . Converges quickly to  $\hat{\theta} = 0.2$ .
- This example is not realistic. Easy analytic solution. Taking the required expectation is trickier in real applications because one needs to know the conditional distribution of the complete data given the missing data.





# 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,  $c = (2a)^{-1}$ ,  $\pi$  is the proportion of the normal distribution in the mixture and  $\theta = (\mu, \sigma, \pi)^T$  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  $\theta$ .



# 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:

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

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

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

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

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

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



# The nature of EM

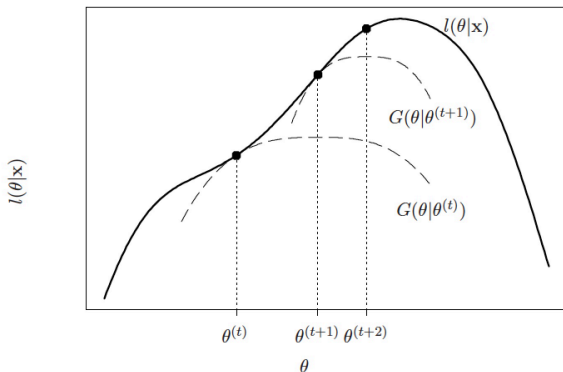
- **Ascent:** Each M-step increases the log likelihood.
- **Convergence:** is linear (slow!). Rate is inversely related to the proportion of missing data.
- **Optimization transfer:**

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

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



# 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

Some variants

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

Variance estimation

- Louis' method
- SEM algorithm





# Overview

EM algorithm

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm



# 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  $f_{\mathbf{Z}|\mathbf{Y}}(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta}^{(t)})$ . 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  $\hat{Q}^{(t+1)}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \frac{1}{m^{(t)}} \sum_{j=1}^{m^{(t)}} \log f_{\mathbf{X}}(\mathbf{X}_j^{(t)}|\boldsymbol{\theta})$ .
- 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

Some variants

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

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  $\theta$ .

- 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**.



# 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  $Q(\theta, \theta^{(t)})$  subject to (or conditional on) a constraint, say

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

where  $\theta^{(t+(s-1)/S)}$  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  $\theta^{(t+1)} = \theta^{(t+S/S)}$  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  $\theta$  into  $S$  subvectors,  $\theta = (\theta_1, \dots, \theta_S)$ .
- In the  $s$ th CM step, maximize  $Q$  with respect to  $\theta_s$  while holding all other components of  $\theta$  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  $\theta$  while holding  $\theta_s$  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.





# 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 convergence.



# Overview

EM algorithm

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm



# Variance of the MLE

- Let  $\hat{\theta}$  be the MLE of  $\theta$ .
- As  $n \rightarrow \infty$ , the limiting distribution of  $\hat{\theta}$  is  $\mathcal{N}(\theta^*, I(\theta^*)^{-1})$ , where  $\theta^*$  is the true value of  $\theta$ , 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 algorithms allows us to estimate  $\hat{\theta}$ , but it does not directly provide an estimate of  $I(\theta^*)$ .
- Direct computation of  $I(\hat{\theta})$  or by  $I_{obs}(\hat{\theta})$  is often difficult.
- Main methods:
  - 1 Louis' method
  - 2 SEM algorithm
  - 3 Bootstrap (to be studied in a later chapter)



# Overview

EM algorithm

Some variants

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

Variance estimation

- Louis' method
- SEM algorithm



# 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 negative both sides, and taking expectations over the conditional distribution of  $\mathbf{X}$  given  $\mathbf{y}$ ,

$$-\ell''(\boldsymbol{\theta}) = \mathbb{E}[-\ell_c''(\boldsymbol{\theta})|\mathbf{y}] - \mathbb{E}\left[-\frac{\partial^2 \log f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} | \mathbf{y}\right]$$

$$\hat{\mathbf{i}}_{\mathbf{Y}}(\boldsymbol{\theta}) = \hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta}) - \hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{Y}}(\boldsymbol{\theta})$$

where

- $\hat{\mathbf{i}}_{\mathbf{Y}}(\boldsymbol{\theta})$  is the observed information,
- $\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta})$  is the conditional expected (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

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



# 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  $\mathbf{x}_i = (\mathbf{y}, \mathbf{z}_i)$ ,  $i = 1, \dots, N$ , where  $\mathbf{y}$  is the observed dataset, and the  $\mathbf{z}_i$  are imputed missing datasets drawn from  $f(\mathbf{z}|\mathbf{y}; \boldsymbol{\theta})$
- Then,

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

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

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





# Overview

EM algorithm

Some variants

Facilitating the E-step

Facilitating the M-step

Variance estimation

Louis' method

SEM algorithm



# Supplemented EM (SEM) algorithm

- Let  $\Psi$  denotes the EM mapping, defined by

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

having fixed point  $\hat{\theta}$  and Jacobian matrix  $\Psi'(\theta)$  with  $(i,j)$ th element equaling  $\frac{\partial \Psi_i(\theta)}{\partial \theta_j}$ .

- It can be shown that

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

- Further use of the missing information principle leads to

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

- SEM is numerically stable and requires little extra work.



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

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

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

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

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



# SEM algorithm

- 1 Run the EM algorithm to convergence, finding  $\hat{\theta}$ .
- 2 Restart the algorithm from some  $\theta^{(0)}$  near to  $\hat{\theta}$ . 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.

