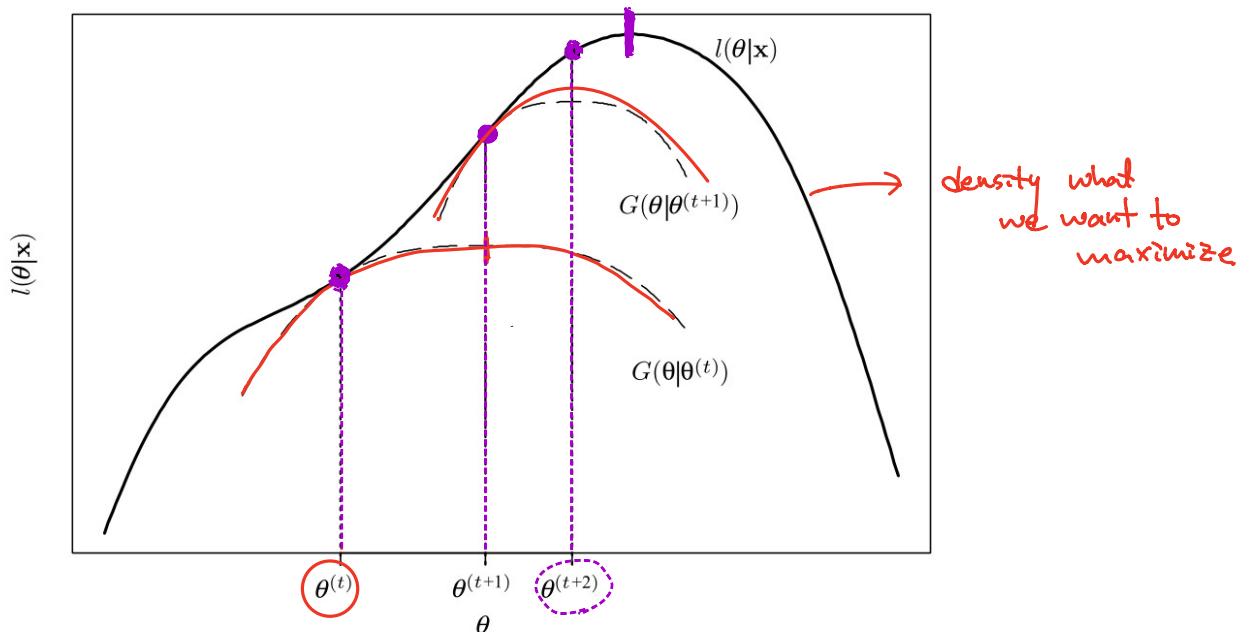


# Appealing Properties of the EM Algorithm

Difficult in derivation  
but easily estimate

- It is typically easily implemented because it relies on complete-data computations
  - The E-step of each iteration only involves taking expectations over complete-data conditional distributions.
  - The M-step of each iteration only requires complete-data maximum likelihood estimation, for which simple closed form expressions are already available.
- It is numerically stable: each iteration is required to increase the log-likelihood  $l_X(\theta)$  in each iteration, and if  $l_X(\theta)$  is bounded, the sequence  $l_X(\theta^{(t)})$  converges to a stationary value.

# One-dimensional Illustration of EM Algorithm



**FIGURE 4.1** One-dimensional illustration of EM algorithm as a minorization or optimization transfer strategy.

# Convergence

The log of the observed-data density can be re-expressed as

$$\log f_X(x | \theta) = \underbrace{\log f_Y(y | \theta)}_{\text{observed}} - \log f_{Z|x}(z | x, \theta).$$

complete      cond z | x

Therefore,

$$E \left\{ \log f_X(x | \theta) | x, \theta^{(t)} \right\} = E \left\{ \log f_Y(y | \theta) | x, \theta^{(t)} \right\} - E \left\{ \log f_{Z|x}(z | x, \theta) | x, \theta^{(t)} \right\},$$

Q( $\theta | \theta^{(t)}$ )      H( $\theta | \theta^{(t)}$ )

where the expectations are taken with respect to the distribution of  $Z | (x, \theta^{(t)})$ . Thus,

$$\log f_X(x | \theta) = Q(\theta | \theta^{(t)}) - H(\theta | \theta^{(t)}),$$

where

$$H(\theta | \theta^{(t)}) = E \left\{ \log f_{Z|x}(Z | x, \theta) | x, \theta^{(t)} \right\}.$$

# Convergence

Show that  $H(\theta | \theta^{(t)})$  is maximized with respect to  $\theta$  when  $\theta = \theta^{(t)}$ . Write

$$H(\theta^{(t)} | \theta^{(t)}) - H(\theta | \theta^{(t)}) \geq 0$$

$$\begin{aligned} H(\theta^{(t)} | \theta^{(t)}) - H(\theta | \theta^{(t)}) &= E \left\{ \log f_{Z|x}(Z | x, \theta^{(t)}) - \log f_{Z|x}(Z | x, \theta) \mid x, \theta^{(t)} \right\} \\ &= \int -\log \left[ \frac{f_{Z|x}(z | x, \theta)}{f_{Z|x}(z | x, \theta^{(t)})} \right] f_{Z|x}(z | x, \theta^{(t)}) dz \end{aligned}$$

$\theta^{(t)} \rightarrow \theta^{(t+1)}$

we know

$$H(\theta^{(t)} | \theta^{(t)}) \text{ maximum} \geq -\log \int f_{Z|x}(z | x, \theta) dz = 0$$

$$H(\theta^{(t+1)} | \theta^{(t)}) < H(\theta^{(t)} | \theta^{(t)})$$

Jensen's inequality

Thus, any  $\theta \neq \theta^{(t)}$  makes  $H(\theta | \theta^{(t)})$  smaller than  $H(\theta^{(t)} | \theta^{(t)})$ .

If we choose  $\theta^{(t+1)}$  to maximize  $Q(\theta | \theta^{(t)})$  with respect to  $\theta$ , then

$$\theta^{(t)} \rightarrow \theta^{(t+1)}$$

$$H(\theta | \theta^{(t)}) \downarrow. \quad \log f_x(x | \theta^{(t+1)}) - \log f_x(x | \theta^{(t)}) \geq 0,$$

$$Q(\theta | \theta^{(t)}) \uparrow$$

since  $Q$  increases and  $H$  decreases, with restrict inequality when

$$Q(\theta^{(t+1)} | \theta^{(t)}) > Q(\theta^{(t)} | \theta^{(t)}).$$

$$\begin{aligned} \log f_x(x | \theta) &= Q(\theta | \theta^{(t)}) - H(\theta | \theta^{(t)}) \\ \therefore f_x(x | \theta^{(t+1)}) &\uparrow \quad \theta^{(t+1)} \downarrow \theta^{(t+2)} \end{aligned}$$

# Convergence

$$\text{portion of missing} = 1 - \text{portion of observed}$$

↑

- Conceptually, the proportion of missing information equals one minus the ratio of the observed information to the information that would be contained in the complete data.
- EM suffers slower convergence when the proportion of missing information is larger.
- The linear convergence of EM can be extremely slow compared to the quadratic convergence of, say, Newton's method, particularly when the fraction of missing information is large.
- However, the ease of implementation and the stable ascent of EM are often very attractive despite its slow convergence.

## Mixed - effect Model

↳ coefficient : fixed effect  
random effect.

$$Y_i = X_i \beta + Z_i b_i + \varepsilon_i \quad b_i, \varepsilon_i \text{ indep.}$$

fixed effect      random effect

(Object: estimate  $\beta$ )

Object: estimate variance of  $b_i$   
 $D$ .

$$b_i \sim N_q(0, D) \quad \varepsilon_i \sim N_{n_i}(0, \sigma^2 I_{n_i})$$



Observed - data log - likelihood

$$l(\beta, D, \sigma^2 | Y_1, \dots, Y_n) = \sum_{i=1}^n \left\{ -\frac{1}{2} (Y_i - X_i \beta)^T \Sigma_i^{-1} (Y_i - X_i \beta) - \frac{1}{2} \log |\Sigma_i| \right\}$$

$$\text{where } \Sigma = Z D Z^T + \sigma^2 I_n.$$

likelihood can be directly maximized for  $(\beta, D, \sigma^2)$

by using Newton-Raphson or Fisher scoring algorithm.

Given  $(D, \sigma^2)$  and hence  $\Sigma_i$ , we obtain  $\beta$  that maximizes the likelihood

by solving.

$$\frac{\partial l(\beta, D, \sigma^2 | Y_1, \dots, Y_n)}{\partial \beta} = \sum_{i=1}^n X_i^T \Sigma_i^{-1} (Y_i - X_i \beta) = 0$$

$$\Rightarrow \hat{\beta} = \left( \sum_{i=1}^n X_i^T \Sigma_i^{-1} X_i \right)^{-1} \sum_{i=1}^n X_i^T \Sigma_i^{-1} Y_i$$

↳ difficult to maximize.  $D, \sigma^2$ .

Complete - data log likelihood.  $\leftarrow$  Need to use EM for  $D, \sigma^2$ .

$\varepsilon_i, b_i$  indep.  $\approx$  Make multivariate normal distribution.

$$- \begin{pmatrix} b_i \\ \varepsilon_i \end{pmatrix} \sim MVN \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} D & 0 \\ 0 & \sigma^2 I_{n_i} \end{pmatrix} \right)$$

$$l_c(\beta, D, \sigma^2 | \varepsilon_1, \dots, \varepsilon_n, b_1, \dots, b_n)$$

$$= \sum_{i=1}^n \left( -\frac{1}{2} b_i^\top D b_i - \frac{1}{2} \log |D| - \frac{1}{2\sigma^2} \varepsilon_i^\top \varepsilon_i - \frac{n}{2} \log \sigma^2 \right)$$

Normal for  $b_i$

Normal for  $\varepsilon_i$

The parameter that maximize the complete-data log-likelihood is obtained as,

conditional on other parameters,

$$\begin{aligned} D &= \frac{1}{N} \sum_{i=1}^n b_i^\top b_i \\ \sigma^2 &= \frac{1}{\sum_{i=1}^n n_i} \sum_{i=1}^n \varepsilon_i^\top \varepsilon_i \\ \beta &= (\sum_{i=1}^n X_i^\top X_i)^{-1} \sum_{i=1}^n X_i^\top (Y_i - Z_i b_i) \end{aligned}$$

$\downarrow Y_i - \text{random effect}$

$\Sigma$ -Step  $E(b_i b_i^\top | Y_i, \beta^{(t)}, D^{(t)}, \sigma^{2(t)})$  Form of  $EX^2$

$E(\varepsilon_i^\top \varepsilon_i | Y_i, \beta^{(t)}, D^{(t)}, \sigma^{2(t)})$

$E(b_i | Y_i, \beta^{(t)}, D^{(t)}, \sigma^{2(t)})$

$\text{Var } X = EX^2 - (EX)^2$

$EX^2 = \underline{\text{Var } X + (EX)^2}$

$$E(b_i b_i^\top | Y_i, \beta^{(t)}, D^{(t)}, \sigma^{2(t)})$$

$$E(b_i b_i^\top | Y_i) = \text{Var}(b_i | Y_i) + E(b_i | Y_i) E(b_i^\top | Y_i)$$

$$\begin{pmatrix} Y_i \\ b_i \end{pmatrix} \sim \text{MVN} \left\{ \begin{pmatrix} X_i \beta \\ 0 \end{pmatrix}, \begin{pmatrix} Z_i D Z_i^\top + \sigma^2 I_n & Z_i D \\ \underline{Z_i^\top} & D \end{pmatrix} \right\}$$

$$\text{Let } \Sigma_i = Z_i D Z_i^\top + \sigma^2 I_n.$$

$$E(b_i | Y_i) = 0 + D Z_i^\top \Sigma_i^{-1} (Y_i - X_i \beta) = \underline{D Z_i^\top \Sigma_i^{-1}} (Y_i - \underline{X_i \beta})$$

$$\text{Var}(b_i | Y_i) = D - D Z_i^\top \Sigma_i^{-1} Z_i D$$

$$E(b_i b_i^\top | Y) = (D - D Z_i^\top \Sigma_i^{-1} Z_i D) + (D Z_i^\top \Sigma_i^{-1} (Y_i - X_i \beta)) (D Z_i^\top \Sigma_i^{-1} (Y_i - X_i \beta))$$

Similarly, we use the relationship

$$E(\varepsilon_i^\top \varepsilon_i | Y_i) = E(\varepsilon_i^\top | Y_i) E(\varepsilon_i | Y_i) + \text{Var}(\varepsilon_i | Y_i)$$

$$\begin{pmatrix} Y_i \\ \varepsilon_i \end{pmatrix} \sim MVN \left( \begin{pmatrix} X_i \beta \\ 0 \end{pmatrix}, \begin{pmatrix} \Sigma & \sigma^2 I_{n_i} \\ \sigma^2 I_{n_i} & \sigma^2 I_{n_i} \end{pmatrix} \right)$$

$$E(\varepsilon_i | Y_i) = 0 + \sigma^2 \Sigma_i^{-1} (Y_i - X_i \beta).$$

$$\text{Var}(\varepsilon_i | Y_i) = \sigma^2 I_{n_i} - \sigma^2 \Sigma_i^{-1}.$$

M-Step  $D^{(t+1)} = \frac{1}{N} \sum_i^n E(b_i b_i^T | Y_i, \beta^{(t)}, D^{(t)}, \sigma^2(t))$

$$\sigma^2(t+1) = \frac{1}{\sum_{i=1}^n n_i} \sum_{i=1}^n E(\varepsilon_i^T \varepsilon_i | Y_i, \beta^{(t)}, D^{(t)}, \sigma^2(t))$$

$$\beta^{(t+1)} = (\sum_{i=1}^n X_i^T X_i)^{-1} \sum_{i=1}^n X_i^T E(Y_i - Z_i b_i | Y_i, \beta^{(t)}, D^{(t)}, \sigma^2(t))$$

### (E) Measurement Error Model

Consider the problem of estimating parameter.

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

$X_i \sim N(\mu_x, \sigma_x^2)$        $\varepsilon_i \sim N(0, \sigma_e^2)$

$\varepsilon_i$  is independent of  $x_i$

In simple regression, we observe  $(x_i, y_i)$

Instead of observing  $(x_i, y_i)$ , suppose we observe  $(z_i, y_i)$ .

in the sample such that  $f(y_i | x_i, z_i) = f(y_i | x)$

Assume that the cond. dist  $g(z_i | x) \sim N(x, \sigma_u^2)$   
given (known)

$$f(x | z, y) \propto f(x | z) f(y | x, z)$$

$$\propto f(x) f(z | x) f(y | x, z)$$

$N(\mu_x, \sigma_x^2)$        $N(x, \sigma_u^2)$        $N(\beta_0 + \beta_1 x, \sigma_e^2)$

$$\propto \exp\left(-\frac{1}{2\sigma_x^2}(x - \mu_x)^2\right) \exp\left(-\frac{1}{2\sigma_u^2}(z - x)^2\right) \exp\left(-\frac{\beta_1^2}{2\sigma_e^2}(x - \beta_1^{-1}y + \beta_0 \beta_1^{-1})^2\right)$$

$$E(x|z, y) = \frac{\mu_x / \sigma_x^2 + z / \sigma_u^2 + \beta_1 (y - \beta_0) / \sigma_e^2}{1 / \sigma_x^2 + 1 / \sigma_u^2 + \beta_1^2 / \sigma_e^2}$$

$$\text{Var}(x|z, y) = \left( \frac{1}{\sigma_x^2} + \frac{1}{\sigma_u^2} + \frac{\beta_1^2}{\sigma_e^2} \right)^{-1}$$

E-Step

$$E(x|z, y; \theta^{(t)}) = \frac{\mu_x^{(t)} / \sigma_x^2(t) + z / \sigma_u^2 + \beta_1^{(t)} (y - \beta_0^{(t)}) / \sigma_e^2(t)}{1 / \sigma_x^2(t) + 1 / \sigma_u^2 + \beta_1^{(t)2} / \sigma_e^2(t)} = \hat{x}^{(t)}$$

$$E(x^2|z, y; \theta^{(t)}) = (\hat{x}^{(t)})^2 + \left( \frac{1}{\sigma_x^2(t)} + \frac{1}{\sigma_u^2} + \beta_1^{(t)2} / \sigma_e^2(t) \right)^{-1}$$

M-Step

$$\mu_x^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \hat{x}^{(t)}$$

$$\sigma_x^2(t+1) = \frac{1}{n} \sum_{i=1}^n E(x^2|z, y; \theta^{(t)}) - (\mu_x^{(t+1)})^2$$

$$\beta_1^{(t+1)} = \frac{1}{\sigma_x^2(t+1)} \left( \frac{1}{n} \sum_{i=1}^n \hat{x}_i^{(t)} y_i - \underbrace{\mu_x^{(t+1)} \hat{\mu}_y}_{= \frac{1}{n} \sum_{i=1}^n y_i} \right)$$

$$\beta_0^{(t+1)} = \hat{\mu}_y - \beta_1^{(t+1)} \mu_x^{(t+1)}$$

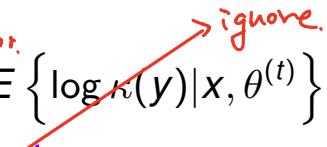
$$\sigma_e^2(t+1) = \hat{\sigma}_y^2 - (\beta_1^{(t+1)})^2 \sigma_x^2(t+1)$$

# Example: Bayesian Posterior Mode

- Consider a Bayesian problem with likelihood  $L(\theta|x)$ , prior  $\pi(\theta)$ , and missing data or parameters  $Z$ . To find the posterior mode, the E-step requires

*Likelihood x prior x marginal density of  $y$*   *normalizing constant.*

$$\begin{aligned} Q(\theta|\theta^{(t)}) &= E \left[ \log \{ L(\theta|y) \pi(\theta) \kappa(y) \} |x, \theta^{(t)} \right] \\ &= E \left\{ \log L(\theta|y) |x, \theta^{(t)} \right\} + \log \pi(\theta) + E \left\{ \log \kappa(y) |x, \theta^{(t)} \right\}, \end{aligned}$$

*complete-data likelihood* 

where the final term is a normalizing constant that can be ignored.

- This function  $Q$  is obtained by simply adding the log prior to the  $Q$  function that would be used in a maximum likelihood setting.
- Unfortunately, the addition of the log prior often makes it more difficult to maximize  $Q$  during the M step.

# Variance Estimation

- In maximum likelihood settings, the EM algorithm is used to find an MLE but does not automatically produce an estimate of the covariance matrix of the MLEs.
- Use the asymptotic normality of the MLEs to justify seeking an estimate of the Fisher information matrix. One way to estimate the covariance matrix is to compute the observed information,  $-I''(\hat{\theta}|x)$ .
- In a Bayesian setting, an estimate of the posterior covariance matrix for  $\theta$  can be motivated by noting the asymptotic normality of the posterior.

# Louis's Method

- Taking second partial derivatives of

$$\log f_X(x | \theta) = Q(\theta | \theta^{(t)}) - H(\theta | \theta^{(t)}),$$

and negating both sides yields

$$-I''(\theta | x) = -Q''(\theta | \omega) \Big|_{\omega=\theta} + H''(\theta | \omega) \Big|_{\omega=\theta},$$
$$\hat{i}_X(\theta) = \hat{i}_Y(\theta) - \hat{i}_{Z|x}(\theta),$$

*observed data* *Fisher info* *complete data* *Fisher info*  *$Z|x$*

where  $\hat{i}_X(\theta) = -I''(\theta | x)$  is the observed information, and  $\hat{i}_Y(\theta)$  and  $\hat{i}_{Z|x}(\theta)$  will be called the complete information and the missing information, respectively.

# Louis's Method

- Interchanging integration and differentiation (when possible), we have

$$\hat{i}_Y(\theta) = -Q''(\theta \mid \omega) \Big|_{\omega=\theta} = -E \{ I''(\theta \mid Y) \mid x, \theta \},$$

which is reminiscent of the Fisher information.

- The missing-information principle can be used to obtain an estimated covariance matrix for  $\hat{\theta}$ . It can be shown that

$$\hat{i}_{Z|x}(\theta) = \text{var} \left\{ \frac{d \log f_{Z|x}(Z \mid x, \theta)}{d\theta} \right\}$$

where the variance is taken with respect to  $f_{Z|x}$ .

# Louis's Method

- Since the expected score is zero at  $\hat{\theta}$ ,

$$\hat{i}_{Z|X}(\hat{\theta}) = \int S_{Z|X}(\hat{\theta}) S_{Z|X}(\hat{\theta})^T f_{Z|X}(z | X, \hat{\theta}) dz,$$

where

$$S_{Z|X}(\theta) = \frac{d \log f_{Z|X}(z | x, \theta)}{d\theta}.$$

- The missing-information principle enables us to express  $\hat{i}_X(\theta)$  in terms of the complete-data likelihood and the conditional density of the missing data given the observed data, while avoiding calculations involving the complicated marginal likelihood of the observed data.
- If  $\hat{i}_Y(\theta)$  or  $\hat{i}_{Z|X}(\theta)$  is difficult to compute analytically, it may be estimated via the Monte Carlo methods.

# SEM Algorithm

- The EM algorithm defines a mapping  $\theta^{(t+1)} = \Psi(\theta^{(t)})$  where the function  $\Psi(\theta) = (\Psi_1(\theta), \dots, \Psi_p(\theta))$  and  $\theta = (\theta_1, \dots, \theta_p)$ . When EM converges, it converges to a fixed point of this mapping, so  $\hat{\theta} = \Psi(\hat{\theta})$  with Jacobian matrix  $\Psi'(\theta)$  with  $(i, j)$ -th element equaling  $d\Psi_i(\theta)/d\theta_j$ . Then,

$$\Psi'(\hat{\theta})^T = \hat{I}_{Z|X}(\hat{\theta}) \hat{I}_Y(\hat{\theta})^{-1}.$$

- If we reexpress the missing information principle as

$$\hat{I}_X(\hat{\theta}) = [I - \hat{I}_{Z|X}(\hat{\theta}) \hat{I}_Y(\hat{\theta})^{-1}] \hat{I}_Y(\hat{\theta}),$$

where  $I$  is an identity matrix, then inverting  $\hat{I}_X(\hat{\theta})$  provides the estimate

$$\hat{\text{var}}\{\hat{\theta}\} = \hat{I}_Y(\hat{\theta})^{-1} \left( I + \Psi'(\hat{\theta})^T [I - \Psi'(\hat{\theta})^T]^{-1} \right).$$

# SEM Algorithm

- 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.
- When coupled with the following numerical differentiation strategy to estimate the increment, Meng and Rubin have termed this approach the supplemented EM (SEM) algorithm.
- Since numerical imprecisions in the differentiation approach affect only the estimated increment, estimation of the covariance matrix is typically more stable than the generic numerical differentiation approach.

# Monte Carlo EM

EM algorithm. [ Improving E-step MCEM.

[ Improving M-Step. ECM / EM Gradient.

In the Monte Carlo EM algorithm, the  $t$ -th E-step can be replaced with the following two steps:

$Y = (X, Z)$   
↳ missing → impute. (multiple imputation)

- ① Draw missing datasets  $Z_1^{(t)}, \dots, Z_{m^{(t)}}^{(t)}$  iid from  $f_{Z|x}(z|x, \theta^{(t)})$ . Each  $Z_j^{(t)}$  is a vector of all the missing values needed to complete the observed dataset, so  $Y_j = (x, Z_j)$  denotes a complete dataset where the missing values have been replaced by  $Z_j$ .
- ② Calculate

$$\hat{Q}^{(t+1)}(\theta|\theta^{(t)}) = \frac{1}{m^{(t)}} \sum_{j=1}^{m^{(t)}} \log f_Y(Y_j^{(t)}|\theta).$$

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

# ECM Algorithm

Improving M-Step

maximization.  
multiple parameter (high-dimension)  
↳ can make cond. density

- Meng and Rubin's ECM algorithm replaces the M step with a series of computationally simpler conditional maximization (CM) steps.
- Each conditional maximization is designed to be a simple optimization problem that constrains  $\theta$  to a particular subspace and permits either an analytical solution or a very elementary numerical solution.
- We call the collection of simpler CM steps after the  $t$ -th E step a CM cycle. Thus, the  $t$ -th iteration of ECM is composed of the  $t$ -th E step and the  $t$ -th CM cycle.

# ECM Algorithm

- ECM Algorithm

- Let  $S$  denote the total number of CM steps in each CM cycle. 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

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

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.
- Clearly any ECM is a GEM algorithm, since each CM step increases  $Q$ .

# ECM Algorithm

- In order for ECM to be convergent, we need to ensure that each CM cycle permits search in any direction for a maximizer of  $Q(\theta|\theta^{(t)})$ , so that ECM effectively maximizes over the original parameter space for  $\theta$  and not over some subspace.
- The art of constructing an effective ECM algorithm lies in choosing the constraints cleverly. Usually, it is natural to partition  $\theta$  into  $S$  sub-vectors,  $\theta = (\theta_1, \dots, \theta_S)$ . Then in the  $s$ th CM step, one might seek to maximize  $Q$  with respect to  $\theta_s$  while holding all other components of  $\theta$  fixed.
- This amounts to the constraint induced by the function  $g_s(\theta) = (\theta_1, \dots, \theta_{s-1}, \theta_{s+1}, \dots, \theta_S)$ . A maximization strategy of this type has previously been termed *iterated conditional modes*.

# EM Gradient Algorithm

- If maximization cannot be accomplished analytically, then one might consider carrying out each M step using an iterative numerical optimization approach.
- This would yield an algorithm that had nested iterative loops.
- The ECM algorithm inserts S conditional maximization steps within each iteration of the EM algorithm, also yielding nested iteration.
- To avoid the computational burden of nested looping, Lange proposed replacing the M step with a single step of Newton's method, thereby approximating the maximum without actually solving for it exactly.

# EM Gradient Algorithm

- The M step is replaced with the update given by

$$\begin{aligned}\theta^{(t+1)} &= \theta^{(t)} - Q''(\theta | \theta^{(t)})^{-1} \left|_{\theta=\theta^{(t)}} Q'(\theta | \theta^{(t)}) \right|_{\theta=\theta^{(t)}} \\ &= \theta^{(t)} - Q''(\theta | \theta^{(t)})^{-1} \left|_{\theta=\theta^{(t)}} I'(\theta^{(t)} | \mathbf{x}) \right|\end{aligned}$$

where  $I'(\theta^{(t)} | \mathbf{x})$  is the evaluation of the score function at the current iterate.

- This EM gradient algorithm has the same rate of convergence to  $\hat{\theta}$  as the full EM algorithm.
- In particular, when  $Y$  has an exponential family distribution with canonical parameter  $\theta$ , ascent is ensured.