

Expectation-Maximization Algorithm

Outline

1. Motivation and EM View
2. The Overview of EM Algorithm
3. Examples
4. Theoretical Issues in EM Algorithm
5. Variants of EM Algorithm

Problem

- Y : p -dimensional random vector with p.d.f. $g(y | \psi)$ on \mathbb{R}^p
- Observe random variable Y modeled by $g(y | \psi)$, $\psi \in \Omega$.
- The likelihood function for ψ formed from the observed data y

$$L(\psi) = g(y | \psi)$$

- Estimate the state of ψ with maximum likelihood.

$$\psi^* = \arg \max_{\psi \in \Omega} \log L(y | \psi)$$

$$\left. \frac{\partial}{\partial \psi} \log L(y | \psi) \right|_{\psi = \psi^*} = 0 \quad (1)$$

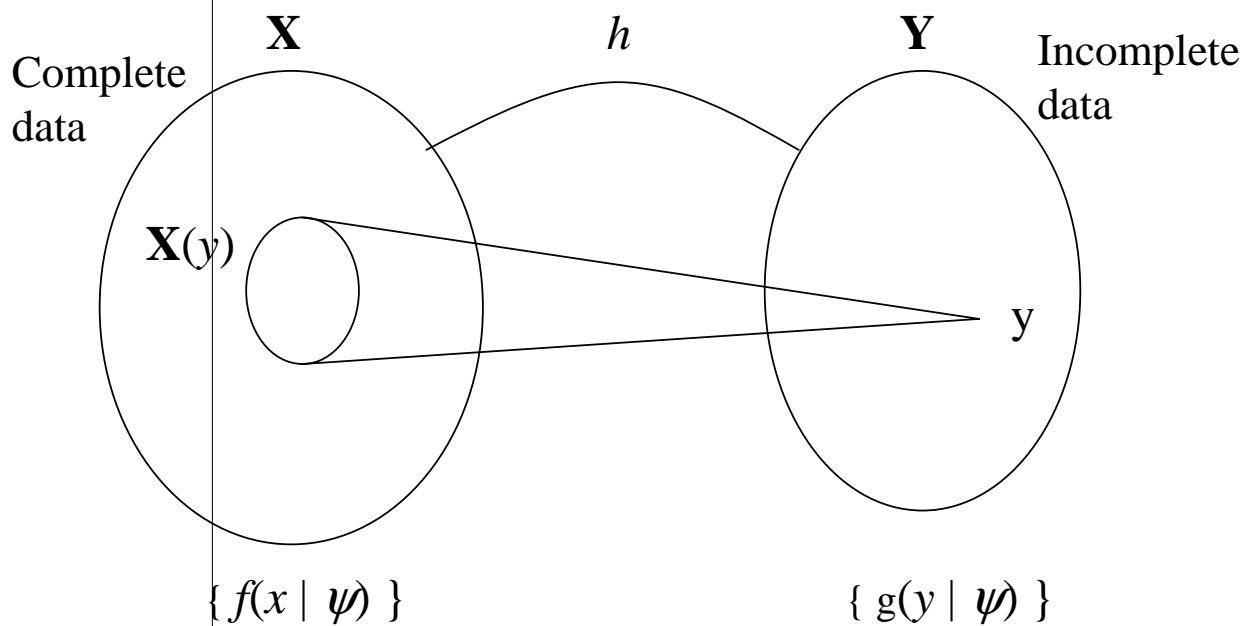
◆ Often it is difficult to solve equation (1).

Motivation

- Data sets with Missing Values
 - ◆ Censored and grouped Observation
 - ◆ Truncated Observation
- Cannot be avoided in many real-world problems

Incomplete Data Model (1)

- Image two sample space \mathbf{X} and \mathbf{Y} , and a map $h: \mathbf{X} \rightarrow \mathbf{Y}$.



Incomplete Data Model (2)

- ◆ Occurrence of $x \in \mathbf{X} \Rightarrow$ occurrence of $y = h(x) \in \mathbf{Y}$.
- ◆ Only y can be observed and $x \in \mathbf{X}(y)$ is revealed.
- ◆ The relation between density f and g is:

$$g(y | \psi) = \int_{\mathbf{X}(y)} f(x | \psi) dx$$

- ◆ In given problems, \mathbf{Y} is fixed while \mathbf{X} can be chosen.

Outline

1. Motivation and EM View
2. The Overview of EM Algorithm
3. Examples
4. Theoretical Issues in EM Algorithm
5. Variants of EM Algorithm

EM Algorithm (1)

- A Standard Tool in the Statistical Repertoire
- Basic Idea
 - ◆ To associate with the given *incomplete-data problem*, a *complete-data problem* for which ML estimation is computationally more tractable
- Y : random vector corresponding to the observed data y , having p.d.f. postulated as $g(y; \psi)$ where ψ is a vector of unknown parameters with parameter space Ω
 - ◆ y : observed incomplete data
 - ◆ x : augmented complete data
 - ◆ z : additional data, referred to as the unobservable or missing data
 - ◆ $f(x; \psi)$: p.d.f. of the random vector X corresponding to the complete-data vector x
 - ◆ Complete-data log likelihood function
$$\log L_c(\psi) = \log f(x; \psi)$$

EM Algorithm (2)

- The expectation of the complete-data log likelihood

$$Q(\psi | \psi') = E[\log f(X | \psi) | y, \psi']$$

- $\psi^{(0)} \in \Omega$: any first approximation to ψ^*
- The EM algorithm consists of repeating two steps.
- On the $(k+1)$ iteration
 - ◆ **E-step.** Calculate $Q(\psi | \psi^{(k)})$.
 - ◆ **M-step.** Choose $\psi^{(k+1)}$ to be any value of $\psi \in \Omega$ that maximizes $Q(\psi, \psi^{(k)})$.

$$\psi^{(k+1)} = \arg \max_{\psi \in \Omega} Q(\psi | \psi^{(k)})$$

Properties

- Choose $\psi^{(k+1)}$ to maximize $\log f(x | \psi)$, where $f(x | \psi)$ is unknown.
 - ◆ Maximize its expectation given the data y and the current fit $\psi^{(k)}$
- Some properties
 - ◆ The constraint $\psi \in \Omega$ is naturally incorporated in the M-step.
 - ◆ The likelihood $g(y | \psi^{(k)})$ is non-decreasing.
 - ◆ Simple and Stable
 - ◆ The sequence of likelihood has a finite limit L^* .
 - L^* can be the local maximum.
- Disadvantage
 - ◆ Slow linear convergence

Outline

1. Motivation and EM View
2. The Overview of EM Algorithm
3. Examples
4. Theoretical Issues in EM Algorithm
5. Variants of EM Algorithm

Multinomial Example (1)

- The observed data vector of frequencies

$$\mathbf{y} = (y_1, y_2, y_3, y_4)^T$$

is postulated to arise from a multinomial distribution with four cells with cell probabilities

$$\frac{1}{2} + \frac{1}{4}\psi, \frac{1}{4}(1-\psi), \frac{1}{4}(1-\psi), \text{ and } \frac{1}{4}\psi$$

with $0 \leq \psi \leq 1$.

- Example

- ◆ $\mathbf{y} = (125, 18, 20, 34)^T$, $n = 197$

Multinomial Example (2)

- The probability function $g(y; \psi)$

$$g(y | \psi) = \frac{n!}{y_1! y_2! y_3! y_4!} \left(\frac{1}{2} + \frac{1}{4}\psi\right)^{y_1} \left(\frac{1}{4} - \frac{1}{4}\psi\right)^{y_2} \left(\frac{1}{4} - \frac{1}{4}\psi\right)^{y_3} \left(\frac{1}{4}\psi\right)^{y_4}$$

- Log likelihood for ψ

$$\log L(\psi) = y_1 \log(2 + \psi) + (y_2 + y_3) \log(1 - \psi) + y_4 \log \psi$$

$$\frac{\partial \log L(\psi)}{\partial \psi} = \frac{y_1}{2 + \psi} - \frac{y_2 + y_3}{1 - \psi} + \frac{y_4}{\psi}$$

- ◆ The likelihood equation can be solved explicitly to find the MLE of ψ .

Multinomial Example (3)

- Suppose that the first of the original four multinomial cells could be split into two subcells having probabilities $\frac{1}{2}$ and $\frac{1}{4}\psi$.

- MLE of ψ on the basis of this split

$$(y_{12} + y_4) / (y_{12} + y_2 + y_3 + y_4)$$

- View y as being incomplete

- Complete-data vector

$$\mathbf{x} = (y_{11}, y_{12}, y_2, y_3, y_4)^T$$

- The cell frequencies in \mathbf{x} :

$$\frac{1}{2}, \frac{1}{4}\psi, \frac{1}{4}(1 - \psi), \frac{1}{4}(1 - \psi), \text{ and } \frac{1}{4}\psi$$

- y_{11} and y_{12} : unobservable and missing data

Multinomial Example (4)

- Over all values of \mathbf{x} such that $y_{11} + y_{12} = y_1$,

$$g(y|\psi) = \sum f(x|\psi)$$

$$f(x|\psi) = C(x) \left(\frac{1}{2}\right)^{y_{11}} \left(\frac{1}{4}\psi\right)^{y_{12}} \left(\frac{1}{4} - \frac{1}{4}\psi\right)^{y_2} \left(\frac{1}{4} - \frac{1}{4}\psi\right)^{y_3} \left(\frac{1}{4}\psi\right)^{y_4}$$

$$C(x) = \frac{n!}{y_{11}! y_{12}! y_2! y_3! y_4!}$$

- Complete-data log likelihood

$$\log L_c(\psi) = (y_{12} + y_4) \log \psi + (y_2 + y_3) \log(1 - \psi)$$

Multinomial Example (5)

- Let $\psi^{(0)}$ be the initial value for ψ .
- E-step:
 - ◆ Filling in for unobservable data by averaging the complete-data log likelihood over conditional distribution given the observed data y
 - ◆ Conditional Expectation of $\log L_c(\psi)$ given y

$$Q(\psi | \psi^{(0)}) = E[\log L_c(\psi) | y, \psi^{(0)}]$$

Multinomial Example (6)

- ◆ Replace y_{11} and y_{12} by their conditional expectations given the observed data \mathbf{y}
 - Y_{11} : random variable corresponding to y_{11}
 - Y_{11} has a binomial distribution with sample size y_1 and probability parameter $\frac{1}{2} / \left(\frac{1}{2} + \frac{1}{4} \psi^{(0)} \right)$
 - Initial conditional expectation of Y_{11} given y_1

$$E_{\psi^{(0)}}(Y_{11} | y_1) = y_{11}^{(0)}$$

$$y_{11}^{(0)} = \frac{1}{2} y_1 / \left(\frac{1}{2} + \frac{1}{4} \psi^{(0)} \right)$$

$$y_{12}^{(0)} = y_1 - y_{11}^{(0)} = \frac{1}{4} y_1 \psi^{(0)} / \left(\frac{1}{2} + \frac{1}{4} \psi^{(0)} \right)$$

Multinomial Example (7)

- M-step
 - ◆ Choose $\psi^{(1)}$ to be the value of ψ that maximizes $Q(\psi; \psi^{(0)})$ w.r.t ψ .

$$\begin{aligned} \psi^{(1)} &= (y_{12}^{(0)} + y_4) / (y_{12}^{(0)} + y_2 + y_3 + y_4) \\ &= (y_{12}^{(0)} + y_4) / (n - y_{11}^{(0)}) \end{aligned}$$

- ◆ In general,

$$\psi^{(k+1)} = (y_{12}^{(k)} + y_4) / (n - y_{11}^{(k)})$$

$$y_{11}^{(k)} = \frac{1}{2} y_1 / \left(\frac{1}{2} + \frac{1}{4} \psi^{(k)} \right)$$

$$y_{12}^{(k)} = y_1 - y_{11}^{(k)}$$

Multinomial Example (8)

Iteration	$\Psi^{(k)}$	$\Psi^{(k)} - \hat{\Psi}$	$r^{(k)}$	$\log L(\Psi^{(k)})$
0	0.500000000	0.126821498	—	64.62974
1	0.608247423	0.018574075	0.1465	67.32017
2	0.624321051	0.002500447	0.1346	67.38292
3	0.626488879	0.000332619	0.1330	67.38408
4	0.626777323	0.000044176	0.1328	67.38410
5	0.626815632	0.000005866	0.1328	67.38410
6	0.626820719	0.000000779	0.1328	67.38410
7	0.626821395	0.000000104	0.1328	67.38410
8	0.626821484	0.000000014	—	67.38410

- ◆ $r(k) = (\Psi^{(k+1)} - \Psi^{(k)}) / (\Psi^{(k)} - \Psi^{(k-1)})$
- ◆ Linear convergence equal to 0.1328

Maximum Likelihood Estimator

- Likelihood $P(D|M)$: probability of data D being generated from model M .
- Learning (induction): inference process of deriving a parameterized model $M=M(w)$ from data D

$$P(M | D) = \frac{P(D | M)P(M)}{P(D)} = P(M) \frac{P(D | M)}{P(D)}$$

- Maximum Likelihood Hypothesis

$$h_{ML} = \arg \max_{h \in H} P(D | h)$$

- Maximum Likelihood Estimator

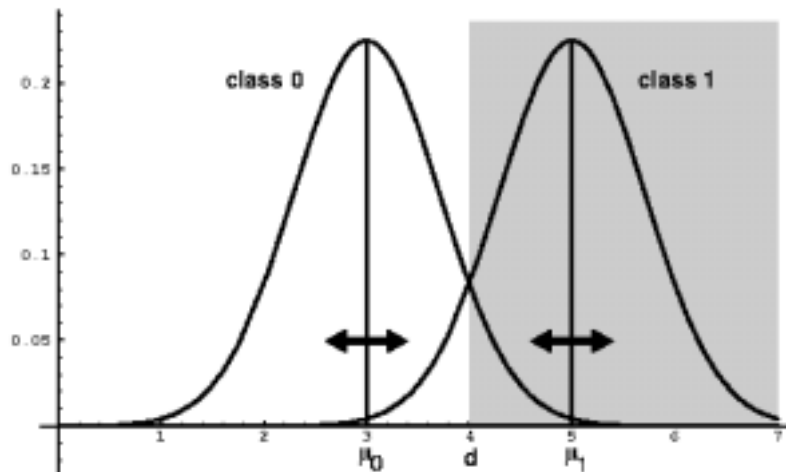
$$\theta^* = \arg \max_{\theta \in H} L(\theta) = \arg \max_{\theta \in H} P(D | \theta)$$

Estimating Mixing Proportions (1)

- Suppose that the p.d.f. of a random vector W has a g -component mixture form

$$f(w|\psi) = \sum_{i=1}^g \pi_i f_i(w)$$

- ◆ $\psi = (\pi_1, \dots, \pi_{g-1})^T$: $g - 1$ mixing proportions



Estimation of Mixing Proportions (2)

- Observed random sample obtained from mixture density

$$\mathbf{y} = (w_1^T, \dots, w_n^T)^T,$$

where

$$w_i = (x, f(x|\psi)) = (x, \sum_{i=1}^g \pi_i f_i(x))$$

- We can observe $f_i(w)$'s, but cannot observe $\pi_i f_i(w)$'s.
 - ◆ We do not know π_i .

Estimation of Mixing Proportions (3)

- Log likelihood function for ψ

$$\begin{aligned}\log L(\psi) &= \sum_{j=1}^n \log f(w_j | \psi) \\ &= \sum_{j=1}^n \log \left\{ \sum_{i=1}^g \pi_i f_i(w_j) \right\}\end{aligned}$$

- ◆ Differentiating with respect to π_i and equating the result to zero

$$\sum_{j=1}^n \left\{ \frac{f_i(w_j)}{f(w_j | \hat{\psi})} - \frac{f_g(w_j)}{f(w_j | \hat{\psi})} \right\} = 0 \quad (i = 1, \dots, g-1)$$

- Does not yield an explicit solution for $\hat{\psi}$

Estimation of Mixing Proportions (4)

- To pose this problem as an incomplete-data one
 - ◆ Introduce as the unobservable or missing data the vector

$$\mathbf{z} = (\mathbf{z}_1^T, \dots, \mathbf{z}_n^T)^T$$

- ◆ \mathbf{z}_i : g -dimensional vector of zero-one indicator variables
- ◆ $z_{ij} = (z_j)_i$: whether w_j arose from the i th component

- MLE of π_I

$$\sum_{j=1}^n z_{ij} / n$$

- Complete-data vector \mathbf{x}

$$\mathbf{x} = (\mathbf{y}^T, \mathbf{z}^T)^T$$

Estimation of Mixing Proportion (5)

- ◆ Complete-data has multinomial form.

$$x_j \sim \text{Mult}_g(1, \{\pi_i f_i(w)\})$$

$$\Pr(\mathbf{X}_j = x_j) = \{\pi_{1i} f_1(w)\}^{z_{1j}} \{\pi_{2i} f_2(w)\}^{z_{2j}} \dots \{\pi_{gi} f_g(w)\}^{z_{gj}}$$

- ◆ Complete-data log likelihood for ψ

$$\log L_c(\psi) = \sum_{i=1}^g \sum_{j=1}^n z_{ij} \log \pi_i + \sum_{i=1}^g \sum_{j=1}^n z_{ij} \log f_i(w_j)$$

Estimation of Mixing Proportions (6)

- E-step

- ◆ Calculation of the current conditional expectation of Z_{ij} given the observed data y

$$\begin{aligned} E(Z_{ij} | y, \psi^{(k)}) &= \text{pr}\{Z_{ij} = 1 | y, \psi^{(k)}\} \\ &= z_{ij}^{(k)} \end{aligned}$$

$$z_{ij}^{(k)} = \frac{\pi_i^{(k)} f_i(w_j)}{f(w_j | \psi^{(k)})} = \frac{\pi_i^{(k)} f_i(w_j)}{\sum_{k=1}^g \pi_k^{(k)} f_k(w_j)}$$

- M-step

$$\pi_i^{(k+1)} = \sum_{j=1}^n z_{ij}^{(k)} / n$$

Outline

1. Motivation and EM View
2. The Overview of EM Algorithm
3. Examples
4. Theoretical Issues in EM Algorithm
5. Variants of EM Algorithm

Non-Decreasing Likelihood (1)

- The conditional density of \mathbf{X} given $\mathbf{Y} = y$

$$k(x | y, \boldsymbol{\psi}) = \frac{f(x | \boldsymbol{\psi})}{\int_{X(y)} f(x | \boldsymbol{\psi}) dx} = \frac{f(x | \boldsymbol{\psi})}{g(y | \boldsymbol{\psi})} \quad x \in X(y)$$

- The log likelihood

$$\begin{aligned} \log L(\boldsymbol{\psi}) &= \log g(y | \boldsymbol{\psi}) \\ &= \log f(x | \boldsymbol{\psi}) - \log k(x | y, \boldsymbol{\psi}) \\ &= \log L_c(\boldsymbol{\psi}) - \log k(x | y, \boldsymbol{\psi}) \end{aligned}$$

- Taking the conditional expectation

$$\begin{aligned} \log L(\boldsymbol{\psi}^{(k)}) &= E[\log L_c(\boldsymbol{\psi}) | y, \boldsymbol{\psi}^{(k)}] - E[\log k(X | y, \boldsymbol{\psi}) | y, \boldsymbol{\psi}^{(k)}] \\ &= Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) - H(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) \end{aligned}$$

Non-Decreasing Likelihood (2)

- What we want show

$$\begin{aligned} & \log L(\psi^{(k+1)}) - \log L(\psi^{(k)}) \\ &= Q(\psi^{(k+1)} | \psi^{(k)}) - Q(\psi^{(k)} | \psi^{(k)}) - \{H(\psi^{(k+1)} | \psi^{(k)}) - H(\psi^{(k)} | \psi^{(k)})\} \\ & \geq 0 \end{aligned}$$

- ◆ The first term

$$Q(\psi^{(k+1)} | \psi^{(k)}) - Q(\psi | \psi^{(k)}) \geq 0$$

- $\psi^{(k+1)}$ is chosen so that for all $\psi \in \Omega$

$$Q(\psi^{(k+1)} | \psi^{(k)}) \geq Q(\psi | \psi^{(k)})$$

- ◆ The second term

$$H(\psi^{(k+1)} | \psi^{(k)}) - H(\psi^{(k)} | \psi^{(k)}) \leq 0 \quad (?)$$

Non-Decreasing Likelihood (3)

- ◆ By Jensen's inequality and the concavity of the logarithmic function

$$\begin{aligned} & H(\psi^{(k+1)} | \psi^{(k)}) - H(\psi^{(k)} | \psi^{(k)}) \\ &= E[\log\{k(X | y, \psi^{(k+1)}) / k(X, y, \psi^{(k)})\} | y, \psi^{(k)}] \\ &\leq \log[E\{k(X | y, \psi^{(k+1)}) / k(X, y, \psi^{(k)})\} | y, \psi^{(k)}] \\ &= \log \int_{X(y)} k(x | y, \psi^{(k+1)}) dx \\ &= 0 \end{aligned}$$

- For a bounded sequence of likelihood value $\{L(\psi^{(k)})\}$, $L(\psi^{(k)})$ converges monotonically to some L^* .

Rate of Convergence (1)

- The EM algorithm implicitly defines a mapping $\psi \rightarrow M(\psi)$

$$\psi^{(k+1)} = M(\psi^{(k)}) \quad (k = 0, 1, 2, \dots)$$

- If $\psi^{(k)}$ converges to some point ψ^* ,

$$\psi^* = M(\psi^*)$$

- In a neighborhood of ψ^* ,

$$\psi^{(k+1)} - \psi^* \approx J(\psi^*)(\psi^{(k)} - \psi^*)$$

where $J(\psi)$ is Jacobian matrix for $M(\psi)$.

- Thus,
$$\frac{\psi^{(k+1)} - \psi^*}{\psi^{(k)} - \psi^*} \approx J(\psi^*)$$

Rate of Convergence (2)

- Near ψ^* , the EM algorithm is a linear iteration with $J(\psi^*)$.

- ◆ $J(\psi^*)$: rate of convergence

- Global rate of convergence

$$r = \lim_{k \rightarrow \infty} \frac{\|\psi^{(k+1)} - \psi^*\|}{\|\psi^{(k)} - \psi^*\|}$$

- ◆ Under certain regularity conditions,

$$r = \lambda_{\max} = \text{the largest eigenvalue of } J(\psi^*)$$

- Global speed of convergence

$$s = 1 - r$$

- ◆ s is the smallest eigenvalue of $S = I_d - J(\psi^*)$

Outline

1. Motivation and EM View
2. The Overview of EM Algorithm
3. Examples
4. Theoretical Issues in EM Algorithm
5. Variants of EM Algorithm

Kinds of EM Variants

1. Generalized EM Algorithm
2. Modification for Maximum *a Posteriori* Estimation
3. ECM Algorithm
4. ECME Algorithm
5. MCEM Algorithm

Generalized EM Algorithm (1)

- M-step of the EM Algorithm

$$\psi^{(k+1)} = \arg \max_{\psi \in \Omega} Q(\psi | \psi^{(k)})$$

- ◆ That is, for all $\psi \in \Omega$,

$$Q(\psi^{(k+1)} | \psi^{(k)}) \geq Q(\psi | \psi^{(k)})$$

- M-step of the GEM Algorithm

- ◆ Choose $\psi^{(k+1)}$ such that

$$Q(\psi^{(k+1)} | \psi^{(k)}) \geq Q(\psi^{(k)} | \psi^{(k)})$$

Generalized EM Algorithm (2)

- Monotonicity of GEM Algorithm

$$\log L(\psi^{(k+1)}) - \log L(\psi^{(k)}) \geq 0$$

- ◆ In the prove of EM Algorithm, only

$$Q(\psi^{(k+1)} | \psi^{(k)}) - Q(\psi | \psi^{(k)}) \geq 0$$

is changed into

$$Q(\psi^{(k+1)} | \psi^{(k)}) - Q(\psi^{(k)} | \psi^{(k)}) \geq 0.$$

For MAP Estimate (1)

- MAP Estimate

$$\psi^* = \arg \max_{\psi \in \Omega} P(\psi | y)$$

$$= \arg \max_{\psi \in \Omega} P(y | \psi)P(\psi)$$

$$\begin{aligned} \therefore \log P(\psi | y) &= \log P(y | \psi) + \log P(\psi) \\ &= \log L(\psi) + \log P(\psi) \end{aligned}$$

For MAP Estimate (2)

- E-step

- ◆ Calculate

$$E[\log p(\psi | x) | y, \psi^{(k)}] = Q(\psi | \psi^{(k)}) + \log p(\psi)$$

- M-step

- ◆ Choose $\psi^{(k+1)}$ to maximize the above equation over $\psi \in \Omega$.

ECM Algorithm (1)

- Motivation
 - ◆ when complete-data ML estimation is rather complicated
 - ◆ Replace a complicated M-step with simpler CM-steps
- CM-steps
 - ◆ Conditional Maximization steps
 - ◆ Maximize the conditional expectation of complete-data log likelihood function found in the preceding E-step subject to constraints on ψ
 - ◆ May require iteration
 - Maximization over smaller dimensional spaces
 - Simpler, faster, and more stable

ECM Algorithm (2)

- M-step is replaced by $S > 1$ steps.
 - ◆ $\psi^{(k+s/S)}$: the value of ψ on the s th CM-step of $(k + 1)$ th iteration
 - ◆ Choose $\psi^{(k+s/S)}$ to maximize $Q(\psi | \psi^{(k)})$ subject to the constraint
$$\mathbf{g}_s(\psi) = \mathbf{g}_s(\psi^{(k+(s-1)/S)})$$
 - ◆ $C = \{\mathbf{g}_s(\psi), s = 1, \dots, S\}$: a set of S preselected (vector) functions
- $\psi^{(k+s/S)}$ satisfies
$$Q(\psi^{(k+s/S)} | \psi^{(k)}) \geq Q(\psi | \psi^{(k)}) \text{ for all } \psi \in \Omega_s(\psi^{(k+(s-1)/S)})$$
where $\Omega_s(\psi^{(k+(s-1)/S)}) \equiv \{\psi \in \Omega : \mathbf{g}_s(\psi) = \mathbf{g}_s(\psi^{(k+(s-1)/S)})\}$
- The value of ψ on the final CM-step
$$\psi^{(k+S/S)} = \psi^{(k+1)}$$

ECM Algorithm (3)

- ECM algorithm is a GEM algorithm.

$$\begin{aligned} Q(\psi^{(k+1)} | \psi^{(k)}) &\geq Q(\psi^{(k+(S-1)/S)} | \psi^{(k)}) \\ &\geq Q(\psi^{(k+(S-2)/S)} | \psi^{(k)}) \\ &\quad \text{⋮} \\ &\geq Q(\psi^{(k)} | \psi^{(k)}) \end{aligned}$$

- Space-Filling Condition

$$\bigcap_{s=1}^S G_s(\psi^{(k)}) = \{\mathbf{0}\} \quad \text{for all } k$$

- ◆ $G_s(\psi)$: column space of $\nabla g_s(\psi)$
- ◆ Complement
 - Convex hull of all feasible direction by $\Omega_s(\psi^{(k+(s-1)/S)})$ is the whole Euclidian space \mathbf{R}^d .
 - Resulting maximization is over the whole parameter space Ω .

ECME Algorithm (4)

- ECME (ECM Either) Algorithm

- ◆ *either* : with this extension, some or all of the CM-steps of the ECM algorithm are replaced by steps that computationally maximize the incomplete-data log likelihood, $\log L(\psi)$, and not the Q -function.

- ◆ Faster convergence.

- CM steps that act on the Q -function must be performed *before* those that act on the actual log likelihood.
- There are situations where the speed of convergence of the ECME algorithm is less than that of the ECM algorithm.

MCEM Algorithm (1)

- When E-step is complex and does not admit a closed-form solution to the computation of $Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$.

- ◆ E-step by a Monte Carlo process

$$\begin{aligned} Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) &= E[\log p(\boldsymbol{\psi} | x) | y] \\ &= \int_{\mathbf{Z}} \log p(\boldsymbol{\psi} | x) p(z | y, \boldsymbol{\psi}^{(k)}) dz \end{aligned}$$

- ◆ \mathbf{Z} : sample space of the missing data
- ◆ $p(z | y, \boldsymbol{\psi})$: conditional density of \mathbf{Z} given y

MCEM Algorithm (2)

- MCEM Algorithm

- ◆ Monte Carlo E-step

- On the k th iteration, draw $z^{(1k)}, \dots, z^{(Mk)}$ from $p(z | y, \boldsymbol{\psi}^{(k)})$

$$Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) = \frac{1}{M} \sum_{m=1}^M \log p(\boldsymbol{\psi} | z^{(mk)}, y)$$

- ◆ M-step

- Maximize $Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$ over $\boldsymbol{\psi}$ to obtain $\boldsymbol{\psi}^{(k+1)}$.

MCEM Algorithm (3)

- An Example

- ◆ Single unknown parameter ψ
- ◆ Single missing variable $z = y_{12}$

$$\text{binomial} \left(y_1, \frac{\frac{1}{4}\psi^{(k)}}{\frac{1}{2} + \frac{1}{4}\psi^{(k)}} \right)$$

- ◆ Complete-data log likelihood

$$\log L_c(\psi) = (y_{12} + y_4) \log \psi + (y_2 + y_3) \log(1 - \psi)$$

- ◆ Q -function in MCEM Algorithm

$$Q(\psi | \psi^{(k)}) = (\bar{z}^{(k)} + y_4) \log \psi + (y_2 + y_3) \log(1 - \psi)$$

$$\psi^{(k+1)} = \frac{\bar{z}^{(k)} + y_4}{\bar{z} + y_4 + y_2 + y_3}$$

MCEM Algorithm (4)

- Difficult problems

- ◆ Monotonicity property is lost.
 - Monitoring of convergence
 - Plot $\psi^{(k)}$ against k .
 - Stabilization of the process indicates the convergence.
 - With fluctuations, the process is terminated or continued with a larger value of M .
- ◆ Choice of M
 - Small values of M in initial stage
 - Increased as the algorithm moves closer to convergence

References

- A. Dempster, N. Laird, and D. Rubin, “Maximum Likelihood from Incomplete Data via the EM algorithm,” *Jour. Roy. Stat. Soc.*, B39, pp.1-38, 1977.
- G. J. McLachlan and T. Krishnan, “*The EM Algorithm and Extensions*,” A Wiley-Interscience Publication, 1997.
- G. McLachlan and D. Peel, “*Finite Mixture Models*,” A Wiley-Interscience Publication, 2000.
- V. Cherkassky and F. Mulier, “*Learning from Data: Concepts, Theory, and Methods*,” A Wiley-Interscience Publication, 1998.