## Lecture Stat 461-561 Expectation-Maximization Algorithm

Arnaud Doucet

January 2008



 $\bullet$  The EM consists of considering maximizing with respect to  $\theta$  the likelihood function

$$L(\theta | \mathbf{y}) = g(\mathbf{y} | \theta) \text{ where } g(\mathbf{y} | \theta) = \int f(\mathbf{x}, \mathbf{y} | \theta) d\mathbf{x}.$$

- We call
  - ${\bf Y}$  the incomplete data (i.e. the observed data) and  $({\bf X}, {\bf Y})$  are the complete data.
- The EM proceeds as

$$\widehat{\boldsymbol{\theta}}^{j+1} = \operatorname*{arg\,max}_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \int \log f\left(\left. \mathbf{x}, \mathbf{y} \right| \boldsymbol{\theta} \right) \ . \ f\left(\left. \mathbf{x} \right| \mathbf{y}, \widehat{\boldsymbol{\theta}}^{j} \right) d\mathbf{x}$$

## Proof of Theorem for Expectation-Maximization Algorithm

- We want to show that  $L\left(\left.\widehat{\theta}^{j+1}\right|\mathbf{y}\right) \geq L\left(\left.\widehat{\theta}^{j}\right|\mathbf{y}\right)$  for  $\widehat{\theta}^{j+1} = \underset{\theta \in \Theta}{\operatorname{arg\,max}} Q\left(\theta, \widehat{\theta}^{j}\right)$ .
- Proof: We have

and for any

$$f(\mathbf{x}|\theta, \mathbf{y}) = \frac{f(\mathbf{x}, \mathbf{y}|\theta)}{g(\mathbf{y}|\theta)} \Leftrightarrow g(\mathbf{y}|\theta) = L(\theta|\mathbf{y}) = \frac{f(\mathbf{x}, \mathbf{y}|\theta)}{f(\mathbf{x}|\theta, \mathbf{y})}$$

thus

$$\log L(\theta | \mathbf{y}) = \log f(\mathbf{x}, \mathbf{y} | \theta) - \log f(\mathbf{x} | \theta, \mathbf{y})$$
  
value  $\hat{\theta}^{j}$ 

$$\log L(\theta | \mathbf{y}) = \underbrace{\int \log f(\mathbf{x}, \mathbf{y} | \theta) . f(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}) d\mathbf{x}}_{=Q(\theta, \widehat{\theta}^{j})} - \int \log f(\mathbf{x} | \theta, \mathbf{y}) . f(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}) d\mathbf{x}$$

• Now the EM ensures by construction that  $Q\left(\widehat{\theta}^{j+1}, \widehat{\theta}^{j}\right) \geq Q\left(\widehat{\theta}^{j}, \widehat{\theta}^{j}\right)$ . So if we can prove that

$$\int \log f\left(\mathbf{x} | \widehat{\theta}^{j+1}, \mathbf{y}\right) . f\left(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}\right) d\mathbf{x} \leq \int \log f\left(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}\right) . f\left(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}\right) d\mathbf{x}$$

then the Theorem is proved.

• We have

$$\int \log \frac{f\left(\mathbf{x} | \widehat{\theta}^{j+1}, \mathbf{y}\right)}{f\left(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}\right)} . f\left(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}\right) d\mathbf{x}$$

$$\leq \log \int \frac{f\left(\mathbf{x} | \widehat{\theta}^{j+1}, \mathbf{y}\right)}{f\left(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}\right)} . f\left(\mathbf{x} | \widehat{\theta}^{j}, \mathbf{y}\right) d\mathbf{x} = 0$$

where we have use Jensen's inequality as log is concave.

• There are numerous variations of the EM in the literature. We will discuss some later.

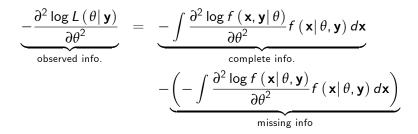
• The EM does not provide an obvious way to compute the observed information matrix given by

$$-\frac{\partial^2 \log L\left(\theta \,|\, \mathbf{y}\right)}{\partial \theta^2}$$

which is an estimate of the inverse covariance matrix of the MLE.

- This quantity is very important in practice and allows us to get some approximate confidence intervals.
- Is it possible to obtain this quantity using EM-type quantities?

• Missing Information Principle



 Proof. It follows straightforwardly from the following identity valid for any x

$$\log L(\theta | \mathbf{y}) = \log f(\mathbf{x}, \mathbf{y} | \theta) - \log f(\mathbf{x} | \theta, \mathbf{y}).$$

• The rate of CV of the EM is highly dependent on the 'ratio' observed info/missing info. The more informative the missing variables are, the slowler the convergence of the algorithm.

## • Proposition : We have

$$\begin{aligned} &-\frac{\partial^2 \log L(\theta | \mathbf{y})}{\partial \theta^2} \\ &= -\int \frac{\partial^2 \log f(\mathbf{x}, \mathbf{y} | \theta)}{\partial \theta^2} f(\mathbf{x} | \theta, \mathbf{y}) d\mathbf{x} - cov\left(\frac{\partial \log f(\mathbf{x}, \mathbf{y} | \theta)}{\partial \theta} \middle| \theta, \mathbf{y}\right) \\ &= -\int \frac{\partial^2 \log f(\mathbf{x}, \mathbf{y} | \theta)}{\partial \theta^2} f(\mathbf{x} | \theta, \mathbf{y}) d\mathbf{x} - \int \frac{\partial \log f(\mathbf{x}, \mathbf{y} | \theta)}{\partial \theta} \frac{\partial \log f(\mathbf{x}, \mathbf{y} | \theta)^{\mathsf{T}}}{\partial \theta} f(\mathbf{x} | \theta, \mathbf{y}) d\mathbf{x} \\ &+ \frac{\partial \log L(\theta | \mathbf{y})}{\partial \theta} \frac{\partial \log L(\theta | \mathbf{y})}{\partial \theta}^{\mathsf{T}} \quad \text{(Louis' identity)} \end{aligned}$$

as

$$\frac{\partial \log L\left(\theta | \mathbf{y}\right)}{\partial \theta} = \int \frac{\partial \log f\left(\mathbf{x}, \mathbf{y} | \theta\right)}{\partial \theta} f\left(\mathbf{x} | \theta, \mathbf{y}\right) d\mathbf{x} \text{ (Fisher's identity)}$$

*Proof.* Fisher's identity is trivial. We have

$$-\int \frac{\partial^2 \log f(\mathbf{x}|\theta, \mathbf{y})}{\partial \theta^2} f(\mathbf{x}|\theta, \mathbf{y}) \, d\mathbf{x} = \int \frac{\partial \log f(\mathbf{x}|\theta, \mathbf{y})}{\partial \theta} \frac{\partial \log f(\mathbf{x}|\theta, \mathbf{y})^{\mathsf{T}}}{\partial \theta} f(\mathbf{x}|\theta, \mathbf{y}) \, d\mathbf{x}$$

and

$$\frac{\partial \log f(\mathbf{x}|\theta, \mathbf{y})}{\partial \theta} = \frac{\partial \log f(\mathbf{x}, \mathbf{y}|\theta)}{\partial \theta} - \frac{\partial \log L(\theta|\mathbf{y})}{\partial \theta}$$

So the result follows directly by noting that  $\int \frac{\partial \log f(\mathbf{x}, \mathbf{y}|\theta)}{\partial \theta} \frac{\partial \log L(\theta|\mathbf{y})^{\mathsf{T}}}{\partial \theta} f(\mathbf{x}|\theta, \mathbf{y}) d\mathbf{x} = \frac{\partial \log L(\theta|\mathbf{y})}{\partial \theta} \frac{\partial \log L(\theta|\mathbf{y})}{\partial \theta}^{\mathsf{T}}.$ 

 Note that Louis' identity is not very friendly.... and is not a direct byproduct of the EM. • Proposition. We have

$$\frac{\partial^{2} \log L\left(\theta \mid \mathbf{y}\right)}{\partial \theta^{2}} = \left. \left\{ \frac{\partial^{2} Q\left(\theta', \theta\right)}{\partial \theta'^{2}} + \frac{\partial^{2} Q\left(\theta', \theta\right)}{\partial \theta' \partial \theta} \right\} \right|_{\theta' = \theta}$$

 $\bullet$  Proof. We show how it is possible to obtain such a result starting from the key identity, for any  $\theta$ 

$$\log L\left(\left.\theta'\right|\mathbf{y}\right) = Q\left(\theta',\theta\right) - \int \log f\left(\mathbf{x}|\left.\theta',\mathbf{y}\right).f\left(\left.\mathbf{x}\right|\theta,\mathbf{y}\right)d\mathbf{x}.$$

Moreover We have

$$\int \frac{\partial \log f(\mathbf{x}|\theta, \mathbf{y})}{\partial \theta} f(\mathbf{x}|\theta, \mathbf{y}) d\mathbf{x} = 0,$$
  
$$\int \frac{\partial^2 \log f(\mathbf{x}|\theta, \mathbf{y})}{\partial \theta^2} f(\mathbf{x}|\theta, \mathbf{y}) d\mathbf{x} = -\int \frac{\partial \log f(\mathbf{x}|\theta, \mathbf{y})}{\partial \theta} \frac{\partial \log f(\mathbf{x}|\theta, \mathbf{y})^{\mathsf{T}}}{\partial \theta} f(\mathbf{x}|\theta, \mathbf{y}) d\mathbf{x}$$
(1)

It follows that

$$\frac{\partial \log L\left(\theta' \mid \mathbf{y}\right)}{\partial \theta'} = \frac{\partial Q\left(\theta', \theta\right)}{\partial \theta'} - \frac{\partial}{\partial \theta'} \int \log f\left(\mathbf{x} \mid \theta', \mathbf{y}\right) \cdot f\left(\mathbf{x} \mid \theta, \mathbf{y}\right) d\mathbf{x}$$
(2)  
thus for  $\theta' = \theta$  we have  $\frac{\partial \log L(\theta \mid \mathbf{y})}{\partial \theta} = \left\{ \frac{\partial Q\left(\theta', \theta\right)}{\partial \theta'} \right\} \Big|_{\theta' = \theta}$ .

• Now differentiating (2) with respect to  $\theta'$  and  $\theta$ , we obtain

$$\frac{\partial \log L\left(\theta' \mid \mathbf{y}\right)}{\partial \theta'^{2}} = \frac{\partial Q\left(\theta', \theta\right)}{\partial \theta'^{2}} - \int \frac{\partial \log f\left(\mathbf{x} \mid \theta', \mathbf{y}\right)}{\partial \theta'^{2}} f\left(\mathbf{x} \mid \theta, \mathbf{y}\right) d\mathbf{x},$$
  
$$\frac{\partial \log L\left(\theta' \mid \mathbf{y}\right)}{\partial \theta' \partial \theta} = 0 = \frac{\partial Q\left(\theta', \theta\right)}{\partial \theta' \partial \theta}$$
  
$$- \int \frac{\partial \log f\left(\mathbf{x} \mid \theta', \mathbf{y}\right)}{\partial \theta'} \cdot \frac{\partial \log f\left(\mathbf{x} \mid \theta, \mathbf{y}\right)^{\mathsf{T}}}{\partial \theta} f\left(\mathbf{x} \mid \theta, \mathbf{y}\right) d\mathbf{x},$$

• Substituting heta= heta', adding the two equations and using (1), we obtain

$$\frac{\partial \log L\left(\theta | \mathbf{y}\right)}{\partial \theta^{2}} = \left. \left\{ \frac{\partial^{2} Q\left(\theta', \theta\right)}{\partial \theta'^{2}} + \frac{\partial^{2} Q\left(\theta', \theta\right)}{\partial \theta' \partial \theta} \right\} \right|_{\theta' = \theta}$$

• This equality is valid at any point  $\theta$ .

• Example: Remember the genetic example where

$$(Y_1, Y_2, Y_3, Y_4) \sim \mathcal{M}\left(n; \frac{1}{2} + \frac{\theta}{4}, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{\theta}{4}\right)$$

• The observed log-likelihood function is given by

$$\log L\left(\theta | \mathbf{y}\right) = cst + y_1 \log \left(2 + \theta\right) + \left(y_2 + y_3\right) \log \left(1 - \theta\right) + y_4 \log \theta.$$

• So we obtain via a direct calculation

$$\frac{\partial \log L\left(\theta \mid \mathbf{y}\right)}{\partial \theta^2} = -\frac{y_1}{\left(2+\theta\right)^2} - \frac{y_2 + y_3}{\left(1-\theta\right)^2} - \frac{y_4}{\theta^2}$$

• Introduce the artificial missing data  $(X_1, X_2)$  such that  $Y_1 = X_1 + X_2$  and define

$$\mathbf{Z} = (X_1, X_2, Y_2, Y_3, Y_4) \sim \mathcal{M}\left(n; \frac{1}{2}, \frac{\theta}{4}, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{\theta}{4}\right).$$

Then

$$\log f\left(\mathbf{z}|\,\theta'\right) = cst + (y_2 + y_3)\log\left(1 - \theta'\right) + (x_2 + y_4)\log\theta'$$
  
and  $\mathbb{E}\left(X_2|\,y_1,\theta\right) = y_1\frac{\theta}{2+\theta}$  so  
 $Q\left(\theta',\theta\right) = cst + (y_2 + y_3)\log\left(1 - \theta'\right) + \left(y_1\frac{\theta}{2+\theta} + y_4\right)\log\theta'$ 

• The second derivatives are given by

$$\frac{\partial^2 Q(\theta',\theta)}{\partial \theta'^2} = -\frac{(y_2+y_3)}{(1-\theta')^2} - \frac{(y_1\frac{\theta}{2+\theta}+y_4)}{\theta'^2},$$
$$\frac{\partial^2 Q(\theta',\theta)}{\partial \theta \partial \theta'} = \frac{2y_1}{(2+\theta)^2}\frac{1}{\theta'}$$

and we can indeed check that 
$$\left\{ \frac{\partial^2 Q\left(\theta',\theta\right)}{\partial \theta'^2} + \frac{\partial^2 Q\left(\theta',\theta\right)}{\partial \theta' \partial \theta} \right\} \Big|_{\theta'=\theta} = \frac{\partial \log L(\theta|\mathbf{y})}{\partial \theta^2}.$$

Arnaud Doucet ()

## The EM as a simple Surrogate Optimization Approach

- The EM approach seems to be closely related to missing data problems... but it can also be seen as a simple surrogate optimization type approach.
- Assume you are interested in maximizing a general function  $f(\theta)$  using an iterative algorithm generating an estimates  $\hat{\theta}^{j}$  at iteration j.
- Assume you can build a function  $g\left( heta,\widehat{ heta}'
  ight)$  such that

$$\begin{array}{rcl} g\left(\theta,\widehat{\theta}^{j}\right) &\leq & f\left(\theta\right) \, \, \text{for any} \, \theta \\ g\left(\widehat{\theta}^{j},\widehat{\theta}^{j}\right) &= & f\left(\widehat{\theta}^{j}\right) \\ \text{then if } \widehat{\theta}^{j+1} = & \text{argmax} \, g\left(\theta,\widehat{\theta}^{j}\right) \, \text{then} \\ & f\left(\widehat{\theta}^{j+1}\right) \geq f\left(\widehat{\theta}^{j}\right). \end{array}$$

• The proof is trivial

$$\begin{split} f\left(\widehat{\theta}^{j+1}\right) - f\left(\widehat{\theta}^{j}\right) \\ &= f\left(\widehat{\theta}^{j+1}\right) - g\left(\widehat{\theta}^{j}, \widehat{\theta}^{j}\right) \\ &= f\left(\widehat{\theta}^{j+1}\right) - g\left(\widehat{\theta}^{j+1}, \widehat{\theta}^{j}\right) + g\left(\widehat{\theta}^{j+1}, \widehat{\theta}^{j}\right) - g\left(\widehat{\theta}^{j}, \widehat{\theta}^{j}\right) \\ &\geq 0 \\ &\text{as } f\left(\widehat{\theta}^{j+1}\right) \geq g\left(\widehat{\theta}^{j+1}, \widehat{\theta}^{j}\right) \text{ and } g\left(\widehat{\theta}^{j+1}, \widehat{\theta}^{j}\right) \geq g\left(\widehat{\theta}^{j}, \widehat{\theta}^{j}\right). \end{split}$$
  
• The EM is a special case where

$$f(\theta) = \log L(\theta | \mathbf{y}),$$
  

$$g(\theta, \hat{\theta}^{j}) = Q(\theta, \hat{\theta}^{j}) + \int \log f(\mathbf{x} | \hat{\theta}^{j}, \mathbf{y}) f(\mathbf{x} | \hat{\theta}^{j}, \mathbf{y}) d\mathbf{x}$$

as

$$\log L(\theta|\mathbf{y}) - \int \log \frac{f(\mathbf{x}|\widehat{\theta},\mathbf{y})}{f(\mathbf{x}|\widehat{\theta}^{j},\mathbf{y})} f(\mathbf{x}|\widehat{\theta}^{j},\mathbf{y}) d\mathbf{x}$$
$$= Q(\theta,\widehat{\theta}^{j}) + \int \log f(\mathbf{x}|\widehat{\theta}^{j},\mathbf{y}) f(\mathbf{x}|\widehat{\theta}^{j},\mathbf{y}) d\mathbf{x}$$

- You have a collection of teams i = 1, ..., N
- Each team *i* plays against the other teams (possibly several times).
- You can only win or lose: no draw.
- We are interesting in ranking the teams.

- We assign to each team *i* a parameter  $\theta_i > 0$ .
- We assume that probability that team *i* beats team *j* is

$$\frac{\theta_i}{\theta_i + \theta_j}$$

• So assuming that this happens  $n_{ij}$  times then the likelihood of  $(\theta_1, ..., \theta_k)$  is

$$\prod_{i,j;i\neq j} \left(\frac{\theta_i}{\theta_i+\theta_j}\right)^{n_{ij}}$$

so for any  $\theta_i^k$ ,  $\theta_j^k$ 

$$I(\theta) = \sum_{i,j;i\neq j} n_{ij} \left(\log \theta_i - \log \left(\theta_i + \theta_j\right)\right).$$

• We use the fact that for any u, v > 0

$$\log \frac{v}{u} \le \frac{v}{u} - 1 \Rightarrow -\log v \ge -\log u - \frac{v - u}{u}$$

so for any  $heta_{i}^{(k)}$  ,  $heta_{j}^{(k)}$ 

$$I(\theta) = \sum_{i,j;i\neq j} n_{ij} \left(\log \theta_i - \log \left(\theta_i + \theta_j\right)\right)$$
  
 
$$\geq \sum_{i,j;i\neq j} n_{ij} \left(\log \theta_i - \log \left(\theta_i^{(k)} + \theta_j^{(k)}\right) - \frac{\left(\theta_i + \theta_j\right) - \left(\theta_i^{(k)} + \theta_j^{(k)} + \theta_j^{(k)}\right)}{\theta_i^{(k)} + \theta_j^{(k)}}\right)$$

• Maximizing the rhs, we obtain

$$\theta_i^{(k+1)} = \frac{\sum_{i \neq j} n_{ij}}{\sum_{i \neq j} (n_{ij} + n_{ji}) / \left(\theta_i^{(k)} + \theta_j^{(k)}\right)}$$

- The key to design this Majorization-Maximization algorithm consists of designing a suitable function g (θ, θ').
- Several 'recipes' are proposed in Hunter&Lange.
- This class of algorithms has been underused in the literature.