# Lecture Notes on the Expectation-Maximization Algorithm

John C. Chao

Econ 721 Lecture Notes

November 22, 2022

November 22, 2022

John C. Chao (Econ 721 Lecture Notes)

 A Motivating Example: Suppose we want to estimate the following two-component mixture model:

Let

$$W_i = \left( {B_i ,\,Y_{1,i} ,\,Y_{2,i} } 
ight)'\,\, {
m for}\,\, i = 1,...,n.$$

Suppose that

$$\{W_i\}\equiv i.i.d.$$

and suppose that

$$\begin{array}{rcl} Y_{1,i} & \sim & \mathcal{N}\left(\mu_{1},\sigma_{1}^{2}\right), & Y_{2,i} \sim \mathcal{N}\left(\mu_{2},\sigma_{2}^{2}\right), \\ B_{i} & = & \left\{ \begin{array}{cc} 1 & \text{with prob } p \\ 0 & \text{with prob } 1-p \end{array} \right. \end{array}$$

Moreover, let

$$Y_i = (1 - B_i) Y_{1,i} + B_i Y_{2,i}$$

and assume that  $B_i$ ,  $Y_{1,i}$ , and  $Y_{2,i}$  are mutually independent and that we only observe  $Y_i$  and not  $B_i$ ,  $Y_{1,i}$ , and  $Y_{2,i}$  separately.

#### Remarks:

- We can think of the underlying generating mechanism is one where a Bernoulli random variable  $B_i \in \{0, 1\}$  is first generated with probability p; and, then, depending on the outcome delivers either  $Y_{1,i}$  or  $Y_{2,i}$ .
- Many economic/econometric models come in the form of a mixture (e.g. models of regime switching, unobserved heterogeneity, etc.)

#### • Probability Density Function:

The pdf of  $Y_i$  is given by

$$g_{Y}\left(y
ight)=\left(1-
ho
ight)\phi_{ heta_{1}}\left(y
ight)+
ho\phi_{ heta_{2}}\left(y
ight)$$
 ,

where for s = 1, 2

$$\begin{array}{rcl} \theta_{s} & = & \left(\mu_{s},\sigma_{s}^{2}\right)', \\ \phi_{\theta_{s}}\left(y\right) & = & \displaystyle\frac{1}{\sigma_{s}\sqrt{2\pi}}\exp\left\{-\displaystyle\frac{1}{2\sigma_{s}^{2}}\left(y-\mu_{s}\right)^{2}\right\}. \end{array}$$

• Additional Notations: Further define

$$\theta = \begin{pmatrix} p \\ \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} p \\ \mu_1 \\ \sigma_1^2 \\ \mu_2 \\ \sigma_2^2 \end{pmatrix} \text{ and } Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}$$

• Log-likelihood Function: Given the notations developed above, the log-likelihood function for this model can be written as

$$\ell\left( heta,Y
ight)=\sum_{i=1}^{n}\ln\left\{\left(1-p
ight)\phi_{ heta_{1}}\left(y
ight)+p\phi_{ heta_{2}}\left(y
ight)
ight\}.$$

Note that this log-likelihood function will be difficult to maximize (say, by Newton's method) because of the sum of terms within the log function.

• The Case Where  $B_i$  is Observed: To understand the idea behind the expectation-maximization algorithm, consider the case where we could observe the values of  $B_i$  (the complete data case). Then, the problem would be a lot easier, since if  $B_i = 1$ ; then,  $Y_i$  comes from model 2; otherwise, it comes from model 1. Hence, if the values of  $B_i$ are observable, then the probability density function of the data would take the form

$$f(B_{i}, Y_{i}|\theta) = f(Y_{i}|B_{i}, \theta_{1}, \theta_{2}) f(B_{i}|p) \\ = \left\{ \left[ \phi_{\theta_{1}}(y_{i}) \right]^{(1-B_{i})} \left[ \phi_{\theta_{2}}(y_{i}) \right]^{B_{i}} \right\} \left\{ [1-p]^{(1-B_{i})} p^{B_{i}} \right\}$$

• Hence, in this case, the likelihood function and the log-likelihood function can be written down, respectively, as

$$L_{0}(\theta, B, Y) = \prod_{i=1}^{n} \left[ \phi_{\theta_{1}}(y_{i}) \right]^{(1-B_{i})} \left[ \phi_{\theta_{2}}(y_{i}) \right]^{B_{i}} [1-p]^{(1-B_{i})} p^{B_{i}},$$
  

$$\ell_{0}(\theta, B, Y) = \sum_{i=1}^{n} \left\{ (1-B_{i}) \ln \phi_{\theta_{1}}(y_{i}) + B_{i} \ln \phi_{\theta_{2}}(y_{i}) \right\}$$
  

$$+ \sum_{i=1}^{n} \left\{ (1-B_{i}) \ln (1-p) + B_{i} \ln p \right\},$$

where  $B = (B_1, ..., B_n)'$ .

• Since in reality the values of  $B_i$  are typically unknown, the EM algorithm proceeds in an iterative manner, substituting for each  $B_i$  with its expected value

$$\begin{split} \gamma_i(\theta) &= E\left[B_i|\theta, Y\right] \\ &= E\left[B_i|\theta, Y_i\right] \text{ (by independence)} \\ &= \Pr\left(B_i = 1|\theta, Y_i\right) \\ &= \frac{f\left(B_i = 1\cap Y_i|\theta\right)}{f\left(Y_i|\theta\right)} \\ &= \frac{f\left(B_i = 1|p\right)f\left(Y_i|B_i = 1, \theta\right)}{f\left(Y_i|\theta\right)} \\ &= \frac{p\phi_{\theta_2}\left(y_i\right)}{\left(1-p\right)\phi_{\theta_1}\left(y_i\right) + p\phi_{\theta_2}\left(y_i\right)}. \end{split}$$

• **Remark:** The quantity  $\gamma_i(\theta) = E[B_i|\theta, Y]$  is often called the **responsibility** of model 2 for observation *i*.

• In addition, taking conditional expectation of the log-likelihood function, we get

$$E\left[\ell_{0}\left(\theta,B,Y\right)|\widehat{\theta},Y\right]$$

$$=\sum_{i=1}^{n}\left\{\left(1-E\left[B_{i}|\widehat{\theta},Y_{i}\right]\right)\ln\phi_{\theta_{1}}\left(y_{i}\right)+E\left[B_{i}|\widehat{\theta},Y_{i}\right]\ln\phi_{\theta_{2}}\left(y_{i}\right)\right\}$$

$$+\sum_{i=1}^{n}\left\{\left(1-E\left[B_{i}|\widehat{\theta},Y_{i}\right]\right)\ln\left(1-p\right)+E\left[B_{i}|\widehat{\theta},Y_{i}\right]\ln p\right\}$$

$$=\sum_{i=1}^{n}\left\{\left(1-\gamma_{i}\left(\widehat{\theta}\right)\right)\ln\phi_{\theta_{1}}\left(y_{i}\right)+\gamma_{i}\left(\widehat{\theta}\right)\ln\phi_{\theta_{2}}\left(y_{i}\right)\right\}$$

$$+\sum_{i=1}^{n}\left\{\left(1-\gamma_{i}\left(\widehat{\theta}\right)\right)\ln\left(1-p\right)+\gamma_{i}\left(\widehat{\theta}\right)\ln p\right\}$$

< 4 ► >

• **Remark:** The EM algorithm iterates back and forth between an expectation step and a maximization step. Under the expectation step, we do a soft assignment of each observation to each model, i.e., the current estimates of the parameters are used to assign responsibilities according to the relative density (under each model) of the sample points. Under the maximization step these responsibilities are used to construct a weighted log-likelihood, which we then maximize to update our estimates of the parameters.

• More precisely, the EM algorithm goes as follows: **Step 1:** Take initial estimates of the parameters  $\hat{\mu}_{1,0}, \hat{\sigma}_{1,0}^2, \hat{\mu}_{2,0}, \hat{\sigma}_{2,0}^2, \hat{p}_0$  (to be specified below). **Step 2:** (Expectation Step) In the  $k^{th}$  step, compute the responsibilities

$$\widehat{\gamma}_{i,k} = \frac{\widehat{p}_{k-1}\phi_{\widehat{\theta}_{2,k-1}}(y_i)}{(1-\widehat{p}_{k-1})\phi_{\widehat{\theta}_{1,k-1}}(y_i) + \widehat{p}_{k-1}\phi_{\widehat{\theta}_{2,k-1}}(y_i)} \text{ for } i = 1, ..., n;$$

where 
$$\widehat{\theta}_{1,k-1} = \left(\widehat{\mu}_{1,k-1}, \widehat{\sigma}_{1,k-1}^2\right)'$$
 and  $\widehat{\theta}_{2,k-1} = \left(\widehat{\mu}_{2,k-1}, \widehat{\sigma}_{2,k-1}^2\right)$ .

• Step 3: (Maximization Step) Compute the weighted means and variances for the  $(k + 1)^{th}$  step as

$$\widehat{\mu}_{1,k} = \frac{\sum_{i=1}^{n} (1 - \widehat{\gamma}_{i,k}) y_i}{\sum_{i=1}^{n} (1 - \widehat{\gamma}_{i,k})}, \widehat{\sigma}_{1,k}^2 = \frac{\sum_{i=1}^{n} (1 - \widehat{\gamma}_{i,k}) (y_i - \widehat{\mu}_{1,k})^2}{\sum_{i=1}^{n} (1 - \widehat{\gamma}_{i,k})}$$

$$\widehat{\mu}_{2,k} = \frac{\sum_{i=1}^{n} \widehat{\gamma}_{i,k} y_i}{\sum_{i=1}^{n} \widehat{\gamma}_{i,k}}, \qquad \widehat{\sigma}_{2,k}^2 = \frac{\sum_{i=1}^{n} (1 - \widehat{\gamma}_{i,k}) (y_i - \widehat{\mu}_{2,k})^2}{\sum_{i=1}^{n} (1 - \widehat{\gamma}_{i,k})}$$

Also, compute the mixing probability as

$$\widehat{p}_{k+1} = \frac{1}{n} \sum_{i=1}^{n} \widehat{\gamma}_{i,k}.$$

**Step 4:** Iterate steps 2 and 3 until convergence.

• **Remark:** Initial estimates  $\hat{\mu}_{1,0}$  and  $\hat{\mu}_{2,0}$  could be made by simply choosing two of the  $y_i$ 's.  $\hat{\sigma}_{1,0}^2$  and  $\hat{\sigma}_{2,0}^2$  could both be set equal to the overall sample variance

$$\frac{1}{n}\sum_{i=1}^{n}\left(y_{i}-\overline{y}\right)^{2}$$

and the initial mixing proportion  $\hat{p}_0$  can be set to 0.5.

 More generally, EM algorithms are useful for problems for which direct maximization of the likelihood function is difficult, but could be made easier by enlarging the sample with latent (unobserved) data. Suppose we have

Y - observed data,

$$heta-$$
 parameter vector,

 $\ell\left( heta, \mathbf{Y}
ight)$  - log-likelihood associated with observed data

and

$$Z -$$
latent or missing data,

$$W = (Y, Z) - \text{ complete data}$$
  
 $\ell_0(\theta; W) - \text{log-likelihood associated with complete data}$ 

- The more general EM algorithm goes as follows:
- Step 1: Take initial estimate of the parameter vector θ<sup>(0)</sup>
   Step 2: (Expectation Step) In the k<sup>th</sup> step, compute the responsibilities

$$Q\left(\theta',\widehat{\theta}^{(k-1)}\right) = E\left[\ell_0\left(\theta',W\right)|Y,\widehat{\theta}^{(k-1)}\right]$$

as a function of the dummy argument  $\theta'$ . **Step 3:** (Maximization Step) Determine  $\hat{\theta}^{(k)}$  as

$$\widehat{ heta}^{(k)} = rg\max_{ heta'} Q\left( heta', \widehat{ heta}^{(k-1)}
ight)$$

**Step 4:** Iterate steps 2 and 3 until convergence.

• Remark: To see why the EM algorithm works, note that

$$f(Z|Y,\theta') = rac{f(Z,Y|\theta')}{f(Y|\theta')} = rac{f(W|\theta')}{f(Y|\theta')},$$

so that

$$f\left(Y| heta'
ight) = rac{f\left(W| heta'
ight)}{f\left(Z|Y, heta'
ight)}$$

In terms of log-likelihoods, this implies that

$$\ell(\theta'; Y) = \ell_0(\theta'; W) - \ln f(Z|Y, \theta')$$
  
=  $\ell_0(\theta'; W) - \ell_1(Y, \theta'|Z)$ 

Taking the conditional expectation with respect to the distribution of W|Y evaluated at the parameter value  $\hat{\theta}$ , we get

$$\ell(\theta'; Y) = E\left[\ell(\theta'; Y) | Y, \widehat{\theta}\right]$$
  
=  $E\left[\ell_0(\theta'; W) | Y, \widehat{\theta}\right] - E\left[\ell_1(Y, \theta' | Z) | Y, \widehat{\theta}\right].$ 

or

•

$$\ell(\theta'; Y) = Q(\theta', \widehat{\theta}) - R(\theta', \widehat{\theta}).$$
Claim:  $R(\theta', \widehat{\theta})$  is maximized at  $R(\widehat{\theta}, \widehat{\theta}).$ 
Proof of Claim: Consider
$$R(\theta', \widehat{\theta}) - R(\widehat{\theta}, \widehat{\theta})$$

$$= E\left[\ell_1(Y, \theta'|Z)|Y, \widehat{\theta}\right] - E\left[\ell_1(Y, \widehat{\theta}|Z)|Y, \widehat{\theta}\right]$$

$$= E\left(\ln\left[\frac{f(Z|Y, \theta')}{f(Z|Y, \widehat{\theta})}\right]|Y, \widehat{\theta}\right)$$

$$\leq \ln\left(E\left[\frac{f(Z|Y, \theta')}{f(Z|Y, \widehat{\theta})}\right]|Y, \widehat{\theta}\right)$$

$$= \ln\left(f(Z|Y, \theta')dZ = \ln 1 = 0.$$

John C. Chao (Econ 721 Lecture Notes)

æ

• Using this result, it follows that if  $\hat{\theta}^{(k)}$  maximizes  $Q\left(\theta', \hat{\theta}^{(k-1)}\right)$ ; then,

$$\ell\left(\theta';Y\right)\Big|_{\theta'=\widehat{\theta}^{(k)}} - \ell\left(\theta';Y\right)\Big|_{\theta'=\widehat{\theta}^{(k-1)}}$$

$$= E\left[\ell\left(\theta';Y\right)|Y,\widehat{\theta}^{(k-1)}\right]\Big|_{\theta'=\widehat{\theta}^{(k)}} - E\left[\ell\left(\theta';Y\right)|Y,\widehat{\theta}^{(k-1)}\right]\Big|_{\theta'=\widehat{\theta}^{(k-1)}}$$

$$= Q\left(\widehat{\theta}^{(k)},\widehat{\theta}^{(k-1)}\right) - Q\left(\widehat{\theta}^{(k-1)},\widehat{\theta}^{(k-1)}\right)$$

$$- \left[R\left(\widehat{\theta}^{(k)},\widehat{\theta}^{(k-1)}\right) - R\left(\widehat{\theta}^{(k-1)},\widehat{\theta}^{(k-1)}\right)\right]$$

$$\geq 0.$$