# 557: MATHEMATICAL STATISTICS II
## THE EM ALGORITHM

The EM Algorithm is a method for producing the maximum likelihood estimates in **incomplete data** problems, that is, models formulated for data that are only partially observed.

Suppose that random variables to be modelled can be partitioned $(\underline{X}, \underline{Y})$ where

- $\underline{X} = (X_1, \ldots, X_m)^{\mathsf{T}}$ are **unobserved**, termed the **augmented data**
- $\underline{X} = (Y_1, \ldots, Y_n)^{\mathsf{T}}$ are **observed**, termed the **incomplete data**
- $(\underline{X}, \underline{Y})$ are termed the **complete data**

where

$$f_{\underline{Y}|\underline{\theta}}(\underline{y}|\underline{\theta}) = \int f_{\underline{X},\underline{Y}|\underline{\theta}}(\underline{x}, \underline{y}|\underline{\theta}) \, d\underline{x}$$

In this formulation,

$$L(\underline{\theta}|\underline{y}) = f_{\underline{Y}|\underline{\theta}}(\underline{y}|\underline{\theta}).$$

is the **incomplete data** likelihood, and

$$L(\underline{\theta}|\underline{x}, \underline{y}) = f_{\underline{X},\underline{Y}|\underline{\theta}}(\underline{x}, \underline{y}|\underline{\theta})$$

is the **complete data** likelihood.

**Algorithm**

The EM Algorithm facilitates maximization of the incomplete data likelihood $L(\underline{\theta}|\underline{y})$ by working with the complete data likelihood $L(\underline{\theta}|\underline{x}, \underline{y})$ and the conditional distribution

$$f_{\underline{X}|\underline{Y},\underline{\theta}}(\underline{x}|\underline{y}, \underline{\theta}) = \frac{f_{\underline{X},\underline{Y}|\underline{\theta}}(\underline{x}, \underline{y}|\underline{\theta})}{f_{\underline{Y}|\underline{\theta}}(\underline{y}|\underline{\theta})} = \frac{L(\underline{\theta}|\underline{x}, \underline{y})}{L(\underline{\theta}|\underline{x})} = K(\underline{x}|\underline{y}, \underline{\theta}) \tag{1}$$

say. It follows from equation (1) that

$$\log L(\underline{\theta}|\underline{y}) = \log L(\underline{\theta}|\underline{x}, \underline{y}) - \log K(\underline{x}|\underline{y}, \underline{\theta}) \tag{2}$$

However, the data $\underline{x}$ are not observed, so consider replacing the right-hand side of equation (2) by the expectations with respect to the conditional density $f_{\underline{X}|\underline{Y},\underline{\theta}}(\underline{x}|\underline{y}, \underline{\theta}')$, for some $\underline{\theta}' \in \Theta$. This yields

$$\log L(\underline{\theta}|\underline{y}) = \mathrm{E}_{f_{\underline{X}|\underline{Y},\underline{\theta}}}[\log L(\underline{\theta}|\underline{X}, \underline{Y})|\underline{y}, \underline{\theta}'] - \mathrm{E}_{f_{\underline{X}|\underline{Y},\underline{\theta}}}[\log K(\underline{X}|\underline{Y}, \underline{\theta})|\underline{y}, \underline{\theta}']. \tag{3}$$

Note that the notation indicates that we condition on a specific (but as yet unspecified) value of $\underline{\theta}'$ when computing the expectations of $\log L(\underline{\theta}|\underline{X}, \underline{y})$ and $\log K(\underline{X}|\underline{y}, \underline{\theta})$ at the $\underline{\theta}$ at which the likelihood on the left-hand side of equation (3) is being computed.

The EM algorithm is an iterative algorithm that produces a sequence of estimates that converges to the (incomplete data) maximum likelihood estimate. Generically, starting from an initial value $\underline{\theta} = \underline{\theta}^{(0)}$, the $(r+1)$st value in the sequence, $\underline{\theta}^{(r+1)}$, is constructed given the $r$th value, $\underline{\theta}^{(r)}$,

$$\underline{\theta}^{(r+1)} = \underset{\underline{\theta} \in \Theta}{\operatorname{argmax}} \, \mathrm{E}_{f_{\underline{X}|\underline{Y},\underline{\theta}}}[\log L(\underline{\theta}|\underline{X}, \underline{Y})|\underline{y}, \underline{\theta}^{(r)}]$$

Two components of this calculation are

- **E-step** : compute the expected conditional log-likelihood
- **M-step** : carry out the maximization of the expectation.

In the traditional notation, we write

$$Q(\underset{\sim}{\theta}|\underset{\sim}{\theta}') = \mathrm{E}_{f_{\underset{\sim}{X}|\underset{\sim}{Y},\underset{\sim}{\theta}}}[\log L(\underset{\sim}{\theta}|\underset{\sim}{X},\underset{\sim}{Y})|\underset{\sim}{y},\underset{\sim}{\theta}']$$

We wish to show that the sequence of estimates produced by

$$\underset{\sim}{\theta}^{(r+1)} = \underset{\underset{\sim}{\theta} \in \Theta}{\mathrm{argmax}}\, Q(\underset{\sim}{\theta}|\underset{\sim}{\theta}^{(r)}) \qquad r = 1, 2, \ldots$$

converges to the maximum likelihood estimate. First, note that for two pdfs $f_1$ and $f_2$ for random variable $Z$, we have by the usual argument that

$$\mathrm{E}_{f_1}[\log f_1(Z)] - \mathrm{E}_{f_1}[\log f_2(Z)] = -\mathrm{E}_{f_1}[\log\{f_2(Z)/f_1(Z)\}] \geq -\log \mathrm{E}_{f_1}[\{f_2(Z)/f_1(Z)\}]$$

$$= -\log \int_{\mathcal{Z}} \{f_2(z)/f_1(z)\} f_1(z)\, dz$$

$$= -\log \int_{\mathcal{Z}} f_2(z)\, dz = 0$$

$$\therefore \qquad \mathrm{E}_{f_1}[\log f_1(Z)] \geq \mathrm{E}_{f_1}[\log f_2(Z)].$$

with equality if and only if $f_1 \equiv f_2$. Hence, for $\underset{\sim}{\theta} \in \Theta$, recalling that

$$K(\underset{\sim}{\theta}|\underset{\sim}{X},\underset{\sim}{Y}) = \frac{L(\underset{\sim}{\theta}|\underset{\sim}{X},\underset{\sim}{Y})}{L(\underset{\sim}{\theta}|\underset{\sim}{Y})} = f_{\underset{\sim}{X}|\underset{\sim}{Y},\underset{\sim}{\theta}}(\underset{\sim}{x}|\underset{\sim}{y},\underset{\sim}{\theta})$$

is itself a (conditional) pdf for all $\underset{\sim}{\theta} \in \Theta$, we have

$$Q(\underset{\sim}{\theta}|\underset{\sim}{\theta}^{(r)}) - \log L(\underset{\sim}{\theta}|\underset{\sim}{y}) = \mathrm{E}_{f_{\underset{\sim}{X}|\underset{\sim}{Y},\underset{\sim}{\theta}}}[\log L(\underset{\sim}{\theta}|\underset{\sim}{X},\underset{\sim}{Y})|\underset{\sim}{y},\underset{\sim}{\theta}^{(r)}] - \log L(\underset{\sim}{\theta}|\underset{\sim}{y})$$

$$= \mathrm{E}_{f_{\underset{\sim}{X}|\underset{\sim}{Y},\underset{\sim}{\theta}}}\left[\log K(\underset{\sim}{\theta}|\underset{\sim}{X},\underset{\sim}{Y})|\underset{\sim}{y},\underset{\sim}{\theta}^{(r)}\right]$$

$$\leq \mathrm{E}_{f_{\underset{\sim}{X}|\underset{\sim}{Y},\underset{\sim}{\theta}}}\left[\log K(\underset{\sim}{\theta}^{(r)}|\underset{\sim}{X},\underset{\sim}{Y})|\underset{\sim}{y},\underset{\sim}{\theta}^{(r)}\right]$$

$$= Q(\underset{\sim}{\theta}^{(r)}|\underset{\sim}{\theta}^{(r)}) - \log L(\underset{\sim}{\theta}^{(r)}|\underset{\sim}{y}).$$

Thus $\log L(\underset{\sim}{\theta}|\underset{\sim}{y}) - Q(\underset{\sim}{\theta}|\underset{\sim}{\theta}^{(r)})$ achieves its **minimum** value when $\underset{\sim}{\theta} = \underset{\sim}{\theta}^{(r)}$. Now suppose that $\underset{\sim}{\theta}^{(r+1)}$ is the value that maximizes $Q(\underset{\sim}{\theta}|\underset{\sim}{\theta}^{(r)})$ over $\Theta$; we have that

$$\log L(\underset{\sim}{\theta}^{(r+1)}|\underset{\sim}{y}) \equiv Q(\underset{\sim}{\theta}^{(r+1)}|\underset{\sim}{\theta}^{(r)}) + \left(\log L(\underset{\sim}{\theta}^{(r+1)}|\underset{\sim}{y}) - Q(\underset{\sim}{\theta}^{(r+1)}|\underset{\sim}{\theta}^{(r)})\right)$$

$$\geq Q(\underset{\sim}{\theta}^{(r)}|\underset{\sim}{\theta}^{(r)}) + \left(\log L(\underset{\sim}{\theta}^{(r)}|\underset{\sim}{y}) - Q(\underset{\sim}{\theta}^{(r)}|\underset{\sim}{\theta}^{(r)})\right)$$

$$= \log L(\underset{\sim}{\theta}^{(r)}|\underset{\sim}{y})$$

and the likelihood attained is **increasing** with the sequence $\underset{\sim}{\theta}^{(0)}, \underset{\sim}{\theta}^{(1)}, \underset{\sim}{\theta}^{(2)}, \ldots$.

**EXAMPLE: Finite Mixture Model**

Suppose that $Y_1 \ldots, Y_n$ are a random sample from the $K$ component finite mixture model

$$f_{Y|\underset{\sim}{\theta}}(y|\underset{\sim}{\theta}) = \sum_{k=1}^{K} \omega_k f_k(y|\theta_k) \qquad y \in \mathbb{R}$$

where $f_1, \ldots, f_K$ are component densities, and

$$0 < \omega_k < 1 \qquad \sum_{k=1}^{K} \omega_k = 1$$

Estimation of $\underset{\sim}{\theta} = (\theta_1, \ldots, \theta_K)^\mathsf{T}$ from the likelihood $L(\underset{\sim}{\theta}|\underset{\sim}{y})$ is in general difficult. However, consider the augmented data $X_1, \ldots, X_n$, where

$$\Pr[X_i = k] = \omega_k \qquad i = 1, \ldots, K$$

are independent random variables so that

$$L(\underset{\sim}{\theta}|\underset{\sim}{X}, \underset{\sim}{Y}) = \prod_{i=1}^{n} \prod_{k=1}^{K} \{\omega_k f_k(y_i|\theta_k)\}^{I_{\{k\}}(X_i)}$$

and

$$\log L(\underset{\sim}{\theta}|\underset{\sim}{X}, \underset{\sim}{Y}) = \sum_{i=1}^{n} \sum_{k=1}^{K} I_{\{k\}}(X_i) \left( \log \omega_k + \log f_k(y_i|\theta_k) \right).$$

The conditional distribution $f_{X|Y,\underset{\sim}{\theta}}(x|y,\underset{\sim}{\theta})$ is a discrete distribution on the set $\{1, 2, \ldots, K\}$ where for each $i = 1, \ldots, n$

$$\Pr[X_i = k|\underset{\sim}{Y}, \underset{\sim}{\omega}, \underset{\sim}{\theta}] = \frac{\omega_k f_k(y_i|\theta_k)}{\sum\limits_{j=1}^{K} \omega_j f_j(y|\theta_j)} = \varpi_k(y_i, \underset{\sim}{\theta}) \qquad k = 1, \ldots, K$$

where $X_1, \ldots, X_n$ are conditionally independent. Thus

$$E_{f_{X_i|Y_i, \underset{\sim}{\theta}, \underset{\sim}{\omega}}}[I_{\{k\}}(X_i)|y_i, \underset{\sim}{\theta}, \underset{\sim}{\omega}] = \varpi_k(y_i, \underset{\sim}{\theta})$$

and hence

$$\begin{aligned}
Q(\underset{\sim}{\theta}, \underset{\sim}{\omega}|\underset{\sim}{\theta}^{(r)}, \underset{\sim}{\omega}^{(r)}) &= E_{f_{\underset{\sim}{X}|\underset{\sim}{Y}, \underset{\sim}{\theta}, \underset{\sim}{\omega}}}[\log L(\underset{\sim}{\theta}|\underset{\sim}{X}, \underset{\sim}{Y})|\underset{\sim}{y}, \underset{\sim}{\theta}^{(r)}, \underset{\sim}{\omega}^{(r)}] \\[2mm]
&= \sum_{i=1}^{n} \sum_{k=1}^{K} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)}) (\log \omega_k + \log f_k(y_i|\theta_k)) \\[2mm]
&= \sum_{k=1}^{K} \left\{ \sum_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)}) \right\} \log \omega_k + \sum_{k=1}^{K} \sum_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)}) \log f_k(y_i|\theta_k) \qquad (4)
\end{aligned}$$

We seek to maximize over $(\underset{\sim}{\theta}, \underset{\sim}{\omega})$ to obtain $(\underset{\sim}{\theta}^{(r+1)}, \underset{\sim}{\omega}^{(r+1)})$ presuming that the values $\varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)})$ are fixed. From the form of equation (4) it is evident that the function is sum of two parts, the first only depending on $\underset{\sim}{\omega}$, the second only dependent on $\underset{\sim}{\theta}$. We can therefore maximize the two parts separately to obtain $(\underset{\sim}{\theta}^{(r+1)}, \underset{\sim}{\omega}^{(r+1)})$.

The first part of equation (4) is of the form of a multinomial likelihood in $\underset{\sim}{\omega}$, therefore, by previous results, it follows that

$$\omega_k^{(r+1)} = \frac{\sum\limits_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)})}{\sum\limits_{j=1}^{K} \sum\limits_{i=1}^{n} \varpi_j^{(r)}(y_i, \underset{\sim}{\theta}^{(r)})} \qquad k = 1, \ldots, K$$

The second part of equation (4) is the sum of $K$ log-likelihoods for the $K$ mixture components which can be maximized separately

$$\theta_k^{(r+1)} = \underset{\theta_k}{\mathrm{argmax}} \sum_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)}) \log f_k(y_i | \theta_k) \tag{5}$$

For certain choices of the component densities, this maximization can be carried out analytically. For example, if $f_k(y|\theta_k)$ is the normal density with expectation $\mu_k$ and variance $\sigma_k^2$, it follows that the new maximizing value equals $\theta_k^{(r+1)} = (\mu_k^{(r+1)}, \sigma_k^{(r+1)})$ where

$$\mu_k^{(r+1)} = \frac{\sum\limits_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)}) y_i}{\sum\limits_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)})}$$

and

$$\sigma_k^{(r+1)} = \sqrt{\frac{\sum\limits_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)})(y_i - \mu_k^{(r+1)})^2}{\sum\limits_{i=1}^{n} \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)})}}$$

Note that in the normal model the terms in (5) correspond to likelihood components of the form

$$\{f_k(y_i|\theta_k)\}^{\varpi_k^{(r)}} = \left(\frac{1}{2\pi\sigma_k^2}\right)^{\varpi_k^{(r)}/2} \exp\left\{-\frac{\varpi_k^{(r)}}{2\sigma_k^2}(y_i - \mu_k^{(r)})^2\right\}$$

so the terms $\varpi_k^{(r)} \equiv \varpi_k^{(r)}(y_i, \underset{\sim}{\theta}^{(r)})$ are acting as weighting factors.