

Lecture Outline

EM Algorithm for MLE (maximum likelihood estimation)

- A review of some theory
- An illustration involving *missing data*
- A result showing that EM's convergence is *monotone*, and giving the *rate of convergence* for the EM algorithm in terms of the amount of *missing information* added.

EM for **MLE** – making a one-step likelihood maximization easier through a (convergent) sequence of simpler maximizations.

Let X_1, X_2, \dots, X_n be *iid* with common density function $p(X | \theta)$.

We are looking to maximize the likelihood function:

$$\hat{\theta} = \operatorname{argmax}_{\Theta} \mathbf{L}(\theta | \mathbf{x}) = \prod_{i=1}^n p(\mathbf{x} | \theta).$$

This may be hard to do as the likelihood function $\mathbf{L}(\theta | \mathbf{x})$ may be complicated.

Instead, it may be easier to work with a likelihood function augmented by data \mathbf{Z}

$$\mathbf{L}(\theta | \mathbf{x}, \mathbf{z})$$

to be integrated out at a later stage of computation.

This is feasible when we can write

$$p(\mathbf{x} | \theta) = \int_{\mathbf{Z}} f(\mathbf{x}, \mathbf{z} | \theta) d\mathbf{z}$$

for some convenient joint density function $f(\mathbf{x}, \mathbf{z} | \theta)$.

Now by the multiplication theorem for densities:

$$f(\mathbf{x}, \mathbf{z} | \theta) = h(\mathbf{z} | \mathbf{x}, \theta) p(\mathbf{x} | \theta)$$

where $h(\mathbf{z} | \mathbf{x}, \theta)$ is a conditional density function for \mathbf{Z} given \mathbf{X} and θ .

It is the convenience of working with the joint density $f(\mathbf{x}, \mathbf{z} | \theta)$ and the conditional density $h(\mathbf{z} | \mathbf{x}, \theta)$ that drives *EM* calculations, as

$$p(\mathbf{x} | \theta) = f(\mathbf{x}, \mathbf{z} | \theta) / h(\mathbf{z} | \mathbf{x}, \theta)$$

Thus, quite generally:

$$(*) \quad \mathbf{\log L}(\theta | \mathbf{x}) = \mathbf{\log L}(\theta | \mathbf{x}, z) - \mathbf{\log h}(z | \mathbf{x}, \theta).$$

Following (Dempster, Laird and Rubin, 1977), with θ_0 arbitrary, define the two functions:

$$(**) \textit{ E-step} \quad Q(\theta | \mathbf{x}, \theta_0) = \int_{\mathbf{Z}} [\mathbf{\log L}(\theta | \mathbf{x}, z)] \mathbf{h}(z | \mathbf{x}, \theta_0) dz$$

and

$$(***) \quad H(\theta | \mathbf{x}, \theta_0) = \int_{\mathbf{Z}} [\mathbf{\log h}(z | \mathbf{x}, \theta_0)] \mathbf{h}(z | \mathbf{x}, \theta_0) dz.$$

$$\mathbf{Then} \quad \mathbf{\log L}(\theta | \mathbf{x}) = Q(\theta | \mathbf{x}, \theta_0) - H(\theta | \mathbf{x}, \theta_0).$$

Begin the iterative process by letting

$$M\text{-step} \quad \hat{\theta}_1 = \operatorname{argmax}_{\Theta} Q(\theta | \mathbf{x}, \theta_0)$$

and then replacing θ_0 with $\hat{\theta}_1$ in (**), which leads to a revised (***) in the light of (*).

Thus,

$$\hat{\theta}_{j+1} = \operatorname{argmax}_{\Theta} Q(\theta | \mathbf{x}, \hat{\theta}_j).$$

(DLR) EM -jargon: $\log \mathbf{L}(\theta | \mathbf{x})$ is the *incomplete* log-likelihood function.

$\log \mathbf{L}(\theta | \mathbf{x}, \mathbf{z})$ is the complete log-likelihood function.

and $Q(\theta | \mathbf{x}, \theta_0)$ is the *expected* log-likelihood function.

Theorem: For the sequence $\hat{\theta}_{j+1} = \operatorname{argmax}_{\Theta} Q(\theta | \mathbf{x}, \hat{\theta}_j)$, $j = 1, \dots$

$$\mathbf{L}(\hat{\theta}_{j+1} | \mathbf{x}) \geq \mathbf{L}(\hat{\theta}_j | \mathbf{x})$$

with equality if and only if $Q(\hat{\theta}_{j+1} | \mathbf{x}, \hat{\theta}_j) = Q(\hat{\theta}_j | \mathbf{x}, \hat{\theta}_j)$.

Proof: Recall that $\log \mathbf{L}(\theta | \mathbf{x}) = Q(\theta | \mathbf{x}, \theta_0) - H(\theta | \mathbf{x}, \theta_0)$.

Then on successive iterations

$$\begin{aligned} \log \mathbf{L}(\hat{\theta}_{j+1} | \mathbf{x}) - \log \mathbf{L}(\hat{\theta}_j | \mathbf{x}) &= \\ & [Q(\hat{\theta}_{j+1} | \mathbf{x}, \hat{\theta}_j) - Q(\hat{\theta}_j | \mathbf{x}, \hat{\theta}_j)] - [H(\hat{\theta}_{j+1} | \mathbf{x}, \hat{\theta}_j) - H(\hat{\theta}_j | \mathbf{x}, \hat{\theta}_j)]. \end{aligned}$$

Evidently $[Q(\hat{\theta}_{j+1} | \mathbf{x}, \hat{\theta}_j) - Q(\hat{\theta}_j | \mathbf{x}, \hat{\theta}_j)] \geq 0$, by the iteration

Thus, we must show that:

$$\int_{\mathbf{Z}} [\mathbf{log} \mathbf{h}(z | \mathbf{x}, \hat{\theta}_{j+1}) - \mathbf{log} \mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)] \mathbf{h}(z | \mathbf{x}, \hat{\theta}_j) dz. \leq 0.$$

Or,

$$\int_{\mathbf{Z}} \mathbf{log} [\mathbf{h}(z | \mathbf{x}, \hat{\theta}_{j+1})/\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)] \mathbf{h}(z | \mathbf{x}, \hat{\theta}_j) dz. \leq 0.$$

Recall, **K-L** information is non-negative and 0 only for identical distributions.

$$\mathbf{E}_{\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)} \mathbf{log} [\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j) / \mathbf{h}(z | \mathbf{x}, \hat{\theta}_{j+1})] \geq 0.$$

Aside: This follows by Jensen's Inequality, twice, noting that for positive rv's $1/\mathbf{E}[\mathbf{X}] < \mathbf{E}[1/\mathbf{X}]$ and that $\mathbf{E}[\mathbf{log} \mathbf{X}] < \mathbf{log} \mathbf{E}[\mathbf{X}]$.

$$\begin{aligned} \text{So, } 0 &\geq -\mathbf{E}_{\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)} \mathbf{log} [\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j) / \mathbf{h}(z | \mathbf{x}, \hat{\theta}_{j+1})] \\ &= \mathbf{E}_{\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)} -\mathbf{log} [\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j) / \mathbf{h}(z | \mathbf{x}, \hat{\theta}_{j+1})] \\ &= \mathbf{E}_{\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)} \mathbf{log} [\mathbf{h}(z | \mathbf{x}, \hat{\theta}_{j+1}) / \mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)] \\ &= \int_{\mathbf{Z}} \mathbf{log} [\mathbf{h}(z | \mathbf{x}, \hat{\theta}_{j+1})/\mathbf{h}(z | \mathbf{x}, \hat{\theta}_j)] \mathbf{h}(z | \mathbf{x}, \hat{\theta}_j) dz \end{aligned}$$

To insure that the sequence $\langle \hat{\theta}_j \rangle$ converges the following result helps:

Theorem: (Boyles, 1983; Wu, 1983)

If the expected log-likelihood function $Q(\theta | \mathbf{x}, \theta_0)$ is continuous in both θ and θ_0 , then all limit points of an *EM* sequence $\langle \hat{\theta}_j \rangle$ are *stationary points* of $\mathbf{L}(\theta | \mathbf{x})$ and $\mathbf{L}(\hat{\theta}_j | \mathbf{x})$ converges monotonically to $\mathbf{L}(\hat{\theta} | \mathbf{x})$ for some *stationary point* $\hat{\theta}$.

That is, then
$$\frac{\partial \log p(\theta | \mathbf{x})}{\partial \theta} \Big|_{\theta = \hat{\theta}} = 0.$$

EM with missing-data.

One-way layout with missing data:

Let X_{ij} denote the response variable of the j^{th} subject among those receiving treatment dose- i .

Statistical model: Assume $X_{ij} \sim N(\mu_i, \sigma^2)$; $i = 1, \dots, k$; $j = 1, \dots, n_i$.

The μ_i are the parameters of interest: average effects of a given treatment dose.

Let $\bar{\mu}$ be an average of average dose effects so that: $\mu_i = \bar{\mu} + \alpha_i$, where $\sum_i \alpha_i = 0$.

That is
$$\bar{\mu} = \sum_i \mu_i / k \quad \text{and} \quad \alpha_i = \mu_i - \bar{\mu}.$$

Note well the relation to the k -MoG problem!

The least squares estimator of μ_i is (evidently) $\bar{x}_i = (1/n_i) \sum_{j=1}^{n_i} x_{ij}$.

And the minimum variance (unbiased) estimators for the other parameters are:

$$\hat{\mu} = (1/k) \sum_i \bar{x}_i \quad \text{and} \quad \hat{\alpha}_i = \bar{x}_i - \hat{\mu}$$

However, when the sample sizes (n_i) are not all equal, the vectors of the coefficients of the X_{ij} in the $\hat{\alpha}_i$ are not orthogonal to the respective vector of coefficients of $\hat{\mu}$. Thus, $\hat{\mu}$ is not independent of the $\hat{\alpha}_i$.

Suppose we have 4 treatment groups, with outcomes

TREATMENTS

T1	T2	T3	T4
x_{11}	x_{21}	x_{31}	x_{41}
x_{12}	x_{22}	x_{32}	x_{42}
z_1	x_{23}	z_3	x_{43}

Observe X_{ij} and use the Z s as the *dummy* missing values to create a balanced sample.

Thus, $X_{ij} \sim N(\bar{\mu} + \alpha_i, \sigma^2)$ and our dimensional parameter $\theta = (\bar{\mu}, \sigma^2, \alpha_1, \alpha_2, \alpha_3, \alpha_4)$.

The *incomplete* likelihood is:

$$\mathbf{L}(\theta | \mathbf{x}) = p(\mathbf{x} | \theta) = \sqrt{(1/2\pi\sigma^2)^{10}} \exp[\sum_{i=1}^4 \sum_{j=1}^{n_i} (x_{ij} - \bar{\mu} - \alpha_i)^2 / \sigma^2]$$

The *complete* likelihood is:

$$\mathbf{L}(\theta | \mathbf{x}, \mathbf{z}) = f(\mathbf{x}, \mathbf{z} | \theta) = \sqrt{(1/2\pi\sigma^2)^{12}} \exp[\sum_{i=1}^4 \sum_{j=1}^3 (x_{ij} - \bar{\mu} - \alpha_i)^2 / \sigma^2]$$

where, of course, $x_{13} = z_1$ and $x_{33} = z_3$.

Now, run the *EM* algorithm with the augmented data (\mathbf{x}, \mathbf{z}) and simplified likelihood (based on a balanced sample) in order to find the MLE for $\mathbf{L}(\theta, \mathbf{x})$.