

Lecture note 13 for STA 250

Shuyang Li

November 13, 2013

Recap: We saw that EM can be used to maximize certain forms of complicated likelihood. The algorithm is as follows:

$$\theta = \arg \max_{\theta} Q(\theta|\theta^{(t)}) \quad (1)$$

$$Q(\theta|\theta^{(t)}) = E[\log [P(y_{obs}, y_{mis}|\theta)] | y_{obs}, \theta^{(t)}] \quad (2)$$

Note: EM maximizes $P(y_{obs}|\theta)$ by using $P(y_{obs}, y_{mis}|\theta)$, which satisfies:

$$\int P(y_{obs}, y_{mis}|\theta) dy_{mis} = P(y_{obs}|\theta) \quad (3)$$

Two key points:

- y_{mis} doesn't have to correspond to "real" missing data;
- choice of y_{mis} is not unique.

One example is as follows:

$$Model : y_{obs}|\theta \sim N(\theta, V + 1) \quad (4)$$

The goal is to find MLE for θ (the answer is y_{obs}). There is no missing data in this example. Consider a "complete" model:

$$y_{obs}|y_{mis} \sim N(y_{mis}, 1) \quad (5)$$

$$y_{mis} \sim N(\theta, V) \quad (6)$$

First we need to check whether $\int P(y_{obs}, y_{mis}|\theta) dy_{mis} = P(y_{obs}|\theta)$. We can show (standard result) that this model fulfils this equation. Here y_{mis} is not "real" missing data. The way to check this equation is to do the expansion as follows:

$$\int P(y_{obs}, y_{mis}|\theta) dy_{mis} = \int P(y_{obs}|y_{mis})P(y_{mis}|\theta) dy_{mis} \quad (7)$$

The model stated here is noted as Model(1). We can also use a different "complete" data model as follows:

$$y_{obs}|y_{\tilde{mis}}, \theta \sim N(y_{\tilde{mis}} + \theta, V) \quad (8)$$

$$y_{\tilde{mis}} \sim N(0, 1) \quad (9)$$

$$(10)$$

we can demonstrate $\int P(y_{obs}, y_{\tilde{mis}}|\theta) dy_{\tilde{mis}} = P(y_{obs}|\theta)$. We denote this model as Model(2). The EM algorithm for Model(1) is:

$$\theta^{t+1} = \frac{1}{V+1}\theta^t + \frac{V}{V+1}y_{obs} \quad (11)$$

The EM algorithm for Model(2) is:

$$\theta^{t+1} = \frac{V}{V+1}\theta^t + \frac{1}{V+1}y_{obs} \quad (12)$$

Hence we have two EMs corresponding to two complete data model, both of which give the same MLE. But which is better?

Model 1 has linear convergence rate $\frac{1}{V+1}$

Model 2 has linear convergence rate $\frac{V}{V+1}$

For the convergence rate, the lower, the better. Hence we choose the one with lower convergence rate as the optimum model. Model 1 is known as a sufficient augmentation scheme as y_{mis} is a sufficient statistics for θ in the "complete" data model. Model 2 is known as an ancillary augmentation scheme as $y_{\tilde{mis}}$ doesn't depend on θ . It turns out that the EM algorithm has an important property: *Monotone convergence*:

$$l(\theta^{t+1}) \geq l(\theta^t) \quad (13)$$

Where

$$l(\theta) = \log P(y_{obs}|\theta)$$

This property makes EM very stable and popular. NR, bisection, scoring don't have this property. Below is the proof of this property.

Note $P(y_{obs}, y_{mis}|\theta) = P(y_{obs}|\theta)P(y_{mis}|y_{obs}, \theta)$. Let $l(\theta) = \log P(y_{obs}|\theta)$, we can get:

$$l(\theta) = \log P(y_{obs}, y_{mis}|\theta) - \log P(y_{mis}|y_{obs}, \theta). \quad (14)$$

Integrate both sides with respect to $P(y_{mis}|y_{obs}, \theta^t)$. Left side = $\int [\log P(y_{obs}|\theta)]P(y_{mis}|y_{obs}, \theta^t)dy_{mis}$. As $\log P(y_{obs}|\theta)$ is not related to y_{mis} ,

$$\int [\log P(y_{obs}|\theta)]P(y_{mis}|y_{obs}, \theta^t)dy_{mis} = \log P(y_{obs}|\theta) = l(\theta) \quad (15)$$

For the right side of equation (14), after integration, the first item in the right side is $Q(\theta|\theta^t)$. Denote the second item after integration as $H(\theta|\theta^t)$, where:

$$H(\theta|\theta^t) = - \int \log P(y_{mis}|y_{obs}, \theta)P(y_{mis}|y_{obs}, \theta^t)dy_{mis} \quad (16)$$

So

$$l(\theta^{t+1}) - l(\theta^t) = [Q(\theta^{t+1}|\theta^t) - Q(\theta^t|\theta^t)] + [H(\theta^{t+1}|\theta^t) - H(\theta^t|\theta^t)] \quad (17)$$

In equation (17), on the right hand side, $\Delta Q = [Q(\theta^{t+1}|\theta^t) - Q(\theta^t|\theta^t)]$ is always ≥ 0 because in the M step, we are maximizing Q . So we only need to show $\Delta H = [H(\theta^{t+1}|\theta^t) - H(\theta^t|\theta^t)] \geq 0$.

Here

$$\Delta H = \int \log \left(\frac{p(y_{mis}|y_{obs}, \theta^t)}{p(y_{mis}|y_{obs}, \theta^{t+1})} \right) P(y_{mis}|y_{obs}, \theta^t)dy_{mis} \quad (18)$$

We know equation (18) is the KL divergence = $\text{KL} [P(y_{mis}|y_{obs}, \theta^t) || P(y_{mis}|y_{obs}, \theta^{t+1})]$. By properties of KL divergence, we know $\Delta H \geq 0$. $\Delta H = 0$ if and only if $P(y_{mis}|y_{obs}, \theta^t) = P(y_{mis}|y_{obs}, \theta^{t+1})$.

Therefore

$$l(\theta^{t+1}) - l(\theta^t) \geq 0 \quad (19)$$

Aside: We can use EM to find posterior models, not just MLE's.

To maximize $\log(\theta|y_{obs})$, Let

$$Q(\theta|y_{obs})_{MAP} = E [P(\theta, y_{mis}|y_{obs})|y_{obs}, \theta^t] \quad (20)$$

$$= \int \log [P(\theta, y_{mis}|y_{obs})] P(y_{mis}|y_{obs}, \theta^t)dy_{mis} \quad (21)$$

Example:

Probit Regression:

$$y_i|x_i, \beta \sim \text{Bin}(1, g(x_i^T \beta)) \quad (22)$$

For logistic regression,

$$g(\mu) = \frac{e^\mu}{1 + e^\mu} \quad (23)$$

For probit regression,

$$g(\mu) = \Phi(\mu) \quad (24)$$

Here Φ is the CDF of $N(0,1)$.

We form a complete data model:

$$y_i|z_i, \beta = I_{z_i \geq 0} \quad (25)$$

$$z_i|\beta \sim N(x_i^T \beta, 1) \quad (26)$$

This model could be connected to the patient's response to some drugs in real application.

Here the complete data is $\{(y_i, z_i), i = 1, \dots, n\}$, the parameter is β , the observed data is $\{y_i, i = 1, \dots, n\}$, the missing data is $\{z_i, i = 1, \dots, n\}$.

First we need to check $\int P(y_i, z_i|\beta)dz_i = P(y_i)$. We need Prof Baines to provide information for this demonstration. Let's derive the EM algorithm for this model:

$$Q(\beta|\beta^t) = E [\log P(y, z|\beta)|y, \beta^t] \quad (27)$$

To take the expectations, we need to know the distribution of $z_i|y_i, \beta^t$.

$$z_i|y_i = 0, \beta^t \sim TN(x_i^T \beta^t, 1, [-\infty, 0]) \quad (28)$$

$$z_i|y_i = 1, \beta^t \sim TN(x_i^T \beta^t, 1, [0, \infty]) \quad (29)$$

$$Q(\beta|\beta^t) = -E \left[\frac{1}{2} \sum_{i=1}^n (z_i - x_i^T \beta)^2 | y_i, \beta^t \right] \quad (30)$$

Maximize $Q(\beta|\beta^t)$.

Let

$$z_i^{t+1} = \begin{cases} x_i^T \beta^t + \frac{\Phi(x_i^T \beta^t)}{1 - \Phi(-x_i^T \beta^t)} & \text{if } y_i = 1 \\ x_i^T \beta^t - \frac{\Phi(x_i^T \beta^t)}{\Phi(-x_i^T \beta^t)} & \text{if } y_i = 0 \end{cases}$$

The maximizer of $Q(\beta|\beta^t)$ with respect to β is seen to be the LS estimator of β when regressing z^{t+1} on x .

$$\beta^{t+1} = (x^T x^{-1}) x^T z^{t+1} \quad (31)$$

where

$$z^{t+1} = (z_1^{t+1}, z_2^{t+1}, \dots, z_n^{t+1})^T \quad (32)$$

- E-step \rightarrow to compute z^{t+1}
- M-step \rightarrow to compute $\beta^{t+1} = (x^T x^{-1}) x^T z^{t+1}$