

MAT-468: Sesión 15, Versiones estocásticas del algoritmo EM

Felipe Osorio

fosorios.mat.utfsm.cl

Departamento de Matemática, UTFSM



Monte Carlo EM (Wei y Tanner, 1990)¹

Para facilitar el paso-E del algoritmo, es posible usar el método Monte Carlo para aproximar la función Q .

Algoritmo 1: Monte Carlo EM (MCEM).

Entrada: Conjunto de datos observados \mathbf{Y}_{obs} y estimación inicial $\boldsymbol{\theta}^{(0)}$.

Salida : Estimación ML de $\boldsymbol{\theta}$.

1 **begin**

2 **Simulación:** Generar $\mathbf{z}_1, \dots, \mathbf{z}_M \stackrel{\text{iid}}{\sim} p(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}; \boldsymbol{\theta}^{(k)})$.

3 **Aproximación:** Sea

$$\widehat{Q}_M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)}) = \frac{1}{M} \sum_{j=1}^M \log p(\boldsymbol{\theta}; \mathbf{z}_j, \mathbf{y}_{\text{obs}}).$$

4 **Paso-M:** actualizar $\boldsymbol{\theta}^{(k+1)}$, como:

$$\boldsymbol{\theta}^{(k+1)} = \arg \max_{\boldsymbol{\theta}} \widehat{Q}_M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)}).$$

5 Iterar entre pasos-E y M hasta alcanzar convergencia.

6 **end**

¹Journal of the American Statistical Association 85, 699-704.



Note que

$$\frac{\partial \widehat{Q}_M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)})}{\partial \boldsymbol{\theta}} = \frac{1}{M} \sum_{j=1}^M \frac{\partial}{\partial \boldsymbol{\theta}} \log p(\boldsymbol{\theta}; \mathbf{z}_j, \mathbf{y}_{\text{obs}})$$
$$\frac{\partial^2 \widehat{Q}_M(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = \frac{1}{M} \sum_{j=1}^M \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log p(\boldsymbol{\theta}; \mathbf{z}_j, \mathbf{y}_{\text{obs}}).$$

De este modo, el **paso-M** del algoritmo MCEM, adopta la forma

$$\boldsymbol{\theta}^{(r+1)} = \boldsymbol{\theta}^{(r)} + \left\{ - \frac{\partial^2 \widehat{Q}_M(\boldsymbol{\theta}^{(r)}; \boldsymbol{\theta}^{(k)})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right\}^{-1} \frac{\partial \widehat{Q}_M(\boldsymbol{\theta}^{(r)}; \boldsymbol{\theta}^{(k)})}{\partial \boldsymbol{\theta}}, \quad (1)$$

a la convergencia de (1) hacemos $\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^*$.

Observación:

Podemos usar los mismos datos simulados $\mathbf{z}_1, \dots, \mathbf{z}_n$ para aproximar

$$\partial Q(\boldsymbol{\theta}; \boldsymbol{\theta}) / \partial \boldsymbol{\theta}, \quad \partial^2 Q(\boldsymbol{\theta}; \boldsymbol{\theta}) / \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top.$$



Sea

$$\mathbf{s}(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) = \nabla_{\boldsymbol{\theta}} \ell_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}),$$

el gradiente de la función de log-verosimilitud de datos completos. Considere la identidad de Fisher,

$$\nabla_{\boldsymbol{\theta}} \ell_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) = \int \mathbf{s}(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) p(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}; \boldsymbol{\theta}) d\mathbf{y}_{\text{mis}}.$$

El método de aproximación estocástica propuesto por Cai (2010) está basado en el algoritmo de Robbins-Monró, usando $\{\gamma_k\}_{k \geq 1}$, tal que

$$\sum_{k=1}^{\infty} \gamma_k = \infty, \quad \sum_{k=1}^{\infty} \gamma_k^2 < \infty.$$

Sea

$$\mathbf{H}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) = -\frac{\partial^2 \ell_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\top}},$$

la matriz de información de datos completos.

²Psychometrika 75, 33-57.

Algoritmo MH-RM (Cai, 2010)

Algoritmo 2: Algoritmo MH-RM.

Entrada: Estimación inicial $\theta^{(0)}$ y Γ_0 matriz definida positiva.

1 **begin**

2 **Simulación:** Generar $z_1, \dots, z_M \stackrel{\text{iid}}{\sim} p(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}; \theta^{(k)})$, y formar el conjunto de datos completos

$$\mathbf{y}_{1,\text{com}}^{(k+1)}, \dots, \mathbf{y}_{M,\text{com}}^{(k+1)}, \quad \mathbf{y}_{i,\text{com}}^{(k+1)} = (z_i, \mathbf{y}_{\text{obs}})$$

3 **Aproximación:** de $\nabla_{\theta} \ell_0(\theta^{(k)}; \mathbf{y}_{\text{obs}})$ y de la matriz de información de datos completos

$$\mathbf{s}_{k+1} = \frac{1}{M} \sum_{j=1}^M \mathbf{s}(\theta^{(k)}; \mathbf{y}_{j,\text{com}}^{(k+1)})$$
$$\Gamma_{k+1} = \Gamma_k + \gamma_k \left\{ \frac{1}{M} \sum_{j=1}^M \mathbf{H}_0(\theta^{(k)}; \mathbf{y}_{\text{com}}^{(k+1)}) - \Gamma_k \right\}.$$

4 **Actualización:** Usamos Robbins-Monr3 para actualizar $\theta^{(k+1)}$,

$$\theta^{(k+1)} = \theta^{(k)} + \gamma_k \Gamma_{k+1}^{-1} \mathbf{s}_{k+1}$$

5 Iterar entre las etapas anteriores hasta alcanzar convergencia.

6 **end**



Algoritmo 3: Algoritmo SAEM.

Entrada: Estimación inicial $\theta^{(0)}$.

1 **begin**

2 **Simulación:** Generar $z_1, \dots, z_M \stackrel{\text{iid}}{\sim} p(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}; \theta^{(k-1)})$, y formar el conjunto de datos completos

$$\mathbf{y}_{1,\text{com}}^{(k+1)}, \dots, \mathbf{y}_{M,\text{com}}^{(k+1)}, \quad \mathbf{y}_{i,\text{com}}^{(k+1)} = (z_i, \mathbf{y}_{\text{obs}})$$

3 **Aproximación:** Actualizar $Q(\theta; \theta^{(k)})$,

$$Q(\theta; \theta^{(k)}) = Q(\theta; \theta^{(k-1)}) + \gamma_k \left\{ \frac{1}{M} \sum_{j=1}^M \log p(\mathbf{y}_{j,\text{com}}^{(k)}; \theta^{(k-1)}) - Q(\theta; \theta^{(k-1)}) \right\}$$

4 **Actualización:** Resolver el problema,

$$\theta^{(k+1)} = \arg \max_{\theta} Q(\theta; \theta^{(k)}).$$

5 Iterar entre las etapas anteriores hasta alcanzar convergencia.

6 **end**

³The Annals of Statistics 27, 94-128.



Algoritmo SAEM (Gu y Kong, 1998)⁴

Algoritmo 4: Algoritmo SAEM.

Entrada: Estimación inicial $\theta^{(0)}$ y Γ_0 matriz definida positiva.

1 **begin**

2 **Simulación:** Generar $z_1, \dots, z_M \stackrel{\text{iid}}{\sim} p(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}; \theta^{(k)})$, y formar el conjunto de datos completos

$$\mathbf{y}_{1,\text{com}}^{(k+1)}, \dots, \mathbf{y}_{M,\text{com}}^{(k+1)}, \quad \mathbf{y}_{i,\text{com}}^{(k+1)} = (z_i, \mathbf{y}_{\text{obs}})$$

3 **Aproximación:** del score y la matriz de información de datos completos

$$\mathbf{g}_{k+1} = \mathbf{g}_k + \gamma_k \left\{ \frac{1}{M} \sum_{j=1}^M \mathbf{s}(\theta^{(k)}; \mathbf{y}_{j,\text{com}}^{(k+1)}) - \mathbf{g}_k \right\}$$

4

$$\Gamma_{k+1} = \Gamma_k + \gamma_k \left\{ \frac{1}{M} \sum_{j=1}^M [\mathbf{H}_c(\theta^{(k)}; \mathbf{y}_{j,\text{com}}^{(k+1)}) - \mathbf{s}(\theta^{(k)}; \mathbf{y}_{j,\text{com}}^{(k+1)}) \mathbf{s}^\top(\theta^{(k)}; \mathbf{y}_{j,\text{com}}^{(k+1)})] - \Gamma_k \right\}.$$

5 **Actualización:** Usar Robbins-Monró para actualizar $\theta^{(k+1)}$,

$$\theta^{(k+1)} = \theta^{(k)} + \gamma_k \Gamma_{k+1}^{-1} \mathbf{g}_{k+1}$$

6 Iterar entre las etapas anteriores hasta alcanzar convergencia.

7 **end**

⁴Proceedings of the National Academy of Sciences of USA 95, 7270-7274.



Disponemos de varias alternativas para estimar el error estándar de $\hat{\theta}$. En particular, podemos usar:

- ▶ Matriz de información de Fisher

$$F_o(\theta) = E\{s_o(\theta)s_o^\top(\theta)\},$$

cuando $\ell_o(\theta)$ es dos veces diferenciable podemos usar

$$F_o(\theta) = E\left\{-\frac{\partial^2 \ell_o(\theta)}{\partial \theta \partial \theta^\top}\right\}.$$

- ▶ Matriz de información observada

$$H_o(\theta) = -\frac{\partial^2 \ell_o(\theta)}{\partial \theta \partial \theta^\top}.$$

- ▶ Versión empírica de la matriz de información

$$\hat{F}_o(\theta) = \frac{1}{n} \sum_{i=1}^n s_i(\theta)s_i^\top(\theta).$$



Principio de información perdida (Orchard y Woodbury, 1972)

Sabemos que

$$\ell_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) = \ell_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) - \log k(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}; \boldsymbol{\theta}),$$

tomando segundas derivadas del negativo de la expresión anterior, tenemos

$$\mathbf{H}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) = \mathbf{H}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) + \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log k(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}; \boldsymbol{\theta}),$$

donde

$$\mathbf{H}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) = -\frac{\partial^2 \ell_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}, \quad \mathbf{H}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) = -\frac{\partial^2 \ell_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top},$$

son las matrices de información (observadas) para el modelo de datos observados y completos, respectivamente.



Principio de información perdida (Orchand y Woodbury, 1972)

Tomando esperanzas con relación a la distribución condicional de $\mathbf{y}_{\text{mis}}|\mathbf{y}_{\text{obs}}$, obtenemos

$$\mathbf{H}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) = \mathbf{F}_c(\boldsymbol{\theta}) - \mathbf{F}_m(\boldsymbol{\theta}), \quad (2)$$

con

$$\mathbf{F}_c(\boldsymbol{\theta}) = \mathbb{E}\{\mathbf{H}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}})|\mathbf{y}_{\text{obs}}\}$$

y

$$\mathbf{F}_m(\boldsymbol{\theta}) = \mathbb{E}\left\{ -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log k(\mathbf{y}_{\text{mis}}|\mathbf{y}_{\text{obs}}; \boldsymbol{\theta}) \middle| \mathbf{y}_{\text{obs}} \right\}.$$

Finalmente, integrando sobre la distribución de \mathbf{y}_{obs} ,

$$\mathbf{F}_o(\boldsymbol{\theta}) = \mathbf{F}_c(\boldsymbol{\theta}) - \mathbb{E}\{\mathbf{F}_m(\boldsymbol{\theta})\}.$$

Note también que, siempre que sea posible intercambiar las operaciones de integración y diferenciación, tenemos

$$\mathbf{F}_c(\boldsymbol{\theta}) = -\frac{\partial^2 Q(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}.$$



Principio de información perdida (Louis, 1982)

Louis (1982) mostró que la matriz de información perdida, puede ser escrita como:

$$\begin{aligned}\mathbf{F}_m(\boldsymbol{\theta}) &= \text{Cov}(\mathbf{s}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}_{\text{obs}}) \\ &= E\{\mathbf{s}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) \mathbf{s}_c^\top(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}_{\text{obs}}\} - E\{\mathbf{s}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}_{\text{obs}}\} E^\top\{\mathbf{s}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}_{\text{obs}}\} \\ &= E\{\mathbf{s}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) \mathbf{s}_c^\top(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}_{\text{obs}}\} - \mathbf{s}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) \mathbf{s}_o^\top(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}})\end{aligned}$$

Substituyendo en Ecuación (2), tenemos

$$\begin{aligned}\mathbf{H}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) &= \mathbf{F}_c(\boldsymbol{\theta}) - \mathbf{F}_m(\boldsymbol{\theta}) \\ &= \mathbf{F}_c(\boldsymbol{\theta}) - E\{\mathbf{s}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) \mathbf{s}_c^\top(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}_{\text{obs}}\} - \mathbf{s}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) \mathbf{s}_o^\top(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}).\end{aligned}$$

A la convergencia del algoritmo EM podemos considerar,

$$\mathbf{H}_o(\boldsymbol{\theta}; \mathbf{y}_{\text{obs}}) = \mathbf{F}_c(\boldsymbol{\theta}) - E\{\mathbf{s}_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) \mathbf{s}_c^\top(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}_{\text{obs}}\}.$$

