

MAT-206: Sesión 8, Estimación: Máxima verosimilitud

Felipe Osorio

fosorios.mat.utfsm.cl

Departamento de Matemática, UTFSM



Definición 1 (Estimador máximo verosímil):

Un estimador $\hat{\theta}_{\text{ML}}$ es llamado estimador máximo verosímil (MLE) de θ , si

$$L(\hat{\theta}_{\text{ML}}; \mathbf{x}) \geq L(\theta; \mathbf{x}), \quad \forall \theta \in \Theta.$$

Es decir, $\hat{\theta}_{\text{ML}}$ debe ser solución del siguiente problema de optimización

$$\max_{\theta \in \Theta} L(\theta; \mathbf{x})$$

o equivalentemente,

$$\max_{\theta \in \Theta} \ell(\theta; \mathbf{x}).$$

Además, en ocasiones escribimos

$$\hat{\theta}_{\text{ML}} := \arg \max_{\theta \in \Theta} \ell(\theta; \mathbf{x}).$$



Resultado 1 (Invarianza del MLE):

Si $\gamma = g(\theta)$ y g es biyectiva. Entonces $\hat{\theta}$ es el MLE para θ si y solo si $\hat{\gamma} = g(\hat{\theta})$ es el MLE para γ .

Demostración:

Considere $L(\theta; x) = f(x; \theta)$ y como g es biyectiva tenemos que

$$\tilde{L}(\gamma; x) = f(x; g^{-1}(\gamma)).$$

Además,

$$\begin{aligned} \tilde{L}(\hat{\gamma}; x) \geq \tilde{L}(\gamma; x), \quad \forall \gamma &\iff f(x; g^{-1}(\hat{\gamma})) \geq f(x; g^{-1}(\gamma)), \quad \forall \gamma \\ \iff f(x; \hat{\theta}) \geq f(x; \theta), \quad \forall \theta &\iff L(\hat{\theta}; x) \geq L(\theta; x), \quad \forall \theta. \end{aligned}$$



Definición 2:

Si $\hat{\theta}_{\text{ML}}$ es el MLE de θ y $\gamma = g(\theta)$. Entonces el MLE de γ es definido como:

$$\hat{\gamma}_{\text{ML}} = g(\hat{\theta}_{\text{ML}}).$$

Observación:

Si $\ell(\theta)$ es continuamente diferenciable, el estimador máximo verosímil $\hat{\theta}_{\text{ML}}$ es dada como una solución de las ecuaciones de verosimilitud:

$$\frac{\partial \ell(\theta)}{\partial \theta} = \mathbf{0},$$

donde $\mathbf{U}(\theta) = \partial \ell(\theta) / \partial \theta$ corresponde al vector score.



Ejemplo:

Considere X_1, \dots, X_n muestra aleatoria desde $N(\mu, \sigma^2)$. De este modo

$$L(\mu, \sigma^2; \mathbf{x}) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\}.$$

Lo que permite obtener la función de log-verosimilitud

$$\ell(\mu, \sigma^2) = -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2,$$

diferenciando con respecto a μ y σ^2 lleva a las ecuaciones de verosimilitud

$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0,$$

$$\frac{\partial \ell(\mu, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0.$$



resolviendo las ecuaciones anteriores para μ y σ^2 , sigue que

$$\hat{\mu}_{\text{ML}} = \bar{x}, \quad \hat{\sigma}_{\text{ML}}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{\text{ML}})^2.$$

TAREA:

Obtener la matriz de información de Fisher de $\theta = (\mu, \sigma^2)^\top$ y luego evaluar en $\theta = \hat{\theta}_{\text{ML}}$.



Ejemplo:

Sea X_1, \dots, X_n una muestra de variables aleatorias IID desde $U[0, \theta]$. Esto es

$$f(x; \theta) = \begin{cases} \frac{1}{\theta}, & 0 \leq x \leq \theta, \\ 0, & \text{en otro caso.} \end{cases} = \frac{1}{\theta} I_{[0, \theta]}(x), \quad \theta > 0.$$

De este modo,¹

$$L(\theta; \mathbf{x}) = \frac{1}{\theta} \prod_{i=1}^n I_{[0, \theta]}(x_i) = \frac{1}{\theta^n} \prod_{i=1}^n I_{[x_i, \infty)}(\theta).$$

Tenemos que

$$\begin{aligned} \prod_{i=1}^n I_{[x_i, \infty)}(\theta) = 1 &\iff I_{[x_i, \infty)}(\theta) = 1, \forall i &\iff x_i \leq \theta, \forall i \\ &\iff \max_i \{x_i\} \leq \theta &\iff I_{[x_{(n)}, \infty)}(\theta) = 1, \end{aligned}$$

donde $x_{(n)} = \max\{x_1, \dots, x_n\}$.

¹Basta notar que $0 \leq x_i \leq \theta \Rightarrow x_i \leq \theta < \infty$.

De este modo, la función

$$L(\theta; \mathbf{x}) = \frac{1}{\theta^n} I_{[x_{(n)}, \infty)}(\theta),$$

es monótona creciente, de ahí que

$$L(\theta; \mathbf{x}) \leq \frac{1}{x_{(n)}^n},$$

y por tanto sigue que $\hat{\theta}_{\text{ML}} = x_{(n)}$.



Ejemplo:

Suponga X_1, \dots, X_n variables aleatorias IID con densidad

$$f(x; a, b) = \frac{b}{2} \exp\{-b|x - a|\}, \quad x \in \mathbb{R},$$

con $a \in \mathbb{R}$ y $b > 0$. De este modo, para b conocido, tenemos

$$\ell(a; \mathbf{x}) = \log(b/2) - b \sum_{i=1}^n |x_i - a|.$$

Es decir, podemos obtener \hat{a}_{ML} , equivalentemente, como la solución de

$$\min_a \sum_{i=1}^n |x_i - a|,$$

y es bien sabido que $\hat{a}_{\text{ML}} = \text{median}\{x_1, \dots, x_n\}$



Método de máxima verosimilitud

Ejemplo:

Sea X_1, \dots, X_n variables IID con distribución común en la FE 1-paramétrica y $\theta \in \Theta$. Considere $\phi = \eta(\theta)$ y sea $\gamma(\phi) = \gamma(\eta(\theta)) = b(\theta)$. Sabemos que

$$L(\phi) = \exp \left[\phi \sum_{i=1}^n T(x_i) - n\gamma(\phi) \right] \prod_{i=1}^n h(x_i).$$

De este modo, la función de log-verosimilitud adopta la forma:

$$\ell(\phi) = \phi \sum_{i=1}^n T(x_i) - n\gamma(\phi) + \sum_{i=1}^n \log h(x_i),$$

lo que lleva a

$$\frac{d\ell(\phi)}{d\phi} = \sum_{i=1}^n T(x_i) - n\gamma'(\phi).$$



Resolviendo la condición de primer orden $d\ell(\phi)/d\phi = 0$, tenemos que $\hat{\phi}_{\text{ML}}$ es solución de la ecuación:

$$\gamma'(\phi) = \frac{1}{n} \sum_{i=1}^n T(x_i).$$

Finalmente por la propiedad de invarianza del MLE sigue que $\hat{\theta}_{\text{ML}} = \eta^{-1}(\hat{\phi}_{\text{ML}})$. Por otro lado, es fácil notar que

$$\frac{d^2 \ell(\phi)}{d\phi^2} = -n\gamma''(\phi) = -nb''(\theta) = -\text{var} \left(\sum_{i=1}^n T(x_i) \right) \leq 0.$$

De lo anterior, sigue que $\ell(\phi)$ es cóncava y por tanto su máximo en $\Phi = \eta(\Theta)$ debe ser único.



Ejemplo:

Suponga X_1, \dots, X_n muestra aleatoria con distribución Weibull, en cuyo caso,

$$f(x; \boldsymbol{\theta}) = \frac{\alpha}{\beta} \left(\frac{x}{\beta}\right)^{\alpha-1} \exp \left\{ - \left(\frac{x}{\beta}\right)^{\alpha} \right\}, \quad x > 0,$$

con $\boldsymbol{\theta} = (\alpha, \beta)^T \in \mathbb{R}_+ \times \mathbb{R}_+$. La función de log-verosimilitud es dada por

$$\ell(\boldsymbol{\theta}) = n(\log \alpha - \log \beta) + (\alpha - 1) \sum_{i=1}^n \log \left(\frac{x_i}{\beta}\right) - \sum_{i=1}^n \left(\frac{x_i}{\beta}\right)^{\alpha}.$$

Diferenciando obtenemos las ecuaciones:

$$\begin{aligned} \frac{n}{\alpha} + \sum_{i=1}^n \log \left(\frac{x_i}{\beta}\right) - \sum_{i=1}^n \left(\frac{x_i}{\beta}\right)^{\alpha} \log \left(\frac{x_i}{\beta}\right) &= 0 \\ -\frac{n\alpha}{\beta} + \frac{\alpha}{\beta} \sum_{i=1}^n \left(\frac{x_i}{\beta}\right)^{\alpha} &= 0. \end{aligned}$$



Método de máxima verosimilitud

Con el objetivo de resolver el problema

$$\max_{\theta \in \Theta} \ell(\theta),$$

suponga la expansión de Taylor en torno de θ^* , como:

$$\ell(\theta^* + \mathbf{p}) = \ell(\theta^*) + \left(\frac{\partial \ell(\theta^*)}{\partial \theta} \right)^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \left(\frac{\partial^2 \ell(\theta^*)}{\partial \theta \partial \theta^\top} \right) \mathbf{p} + o(\|\mathbf{p}\|^2),$$

donde $o(u)$ es un término de error de orden menor que u , conforme $u \rightarrow 0$, es decir,

$$\lim_{u \rightarrow 0} \frac{o(u)}{u} = 0.$$

Defina la función cuadrática,

$$q_k(\mathbf{p}) = \ell(\theta^{(k)}) + \mathbf{U}^\top(\theta^{(k)}) \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \mathbf{H}(\theta^{(k)}) \mathbf{p},$$

donde $\mathbf{U}(\theta) = \partial \ell(\theta) / \partial \theta$ y $\mathbf{H}(\theta) = \partial^2 \ell(\theta) / \partial \theta \partial \theta^\top$.



Método de máxima verosimilitud

Minimizando $q_k(\mathbf{p})$ con relación a \mathbf{p} , lleva al sistema de ecuaciones $\partial q_k(\mathbf{p})/\partial \mathbf{p} = \mathbf{0}$. Es decir, obtenemos:

$$\mathbf{H}(\boldsymbol{\theta}^{(k)})\mathbf{p} = -\mathbf{U}(\boldsymbol{\theta}^{(k)}). \quad (1)$$

Métodos **tipo-Newton** adoptan la forma:

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + \lambda_k \mathbf{p}_k, \quad k = 0, 1, \dots,$$

donde \mathbf{p}_k es la dirección de búsqueda dada por la solución del sistema dado en (1), mientras que λ_k es un largo de paso que debe ser escogido para garantizar que

$$\ell(\boldsymbol{\theta}^{(k)} + \lambda_k \mathbf{p}_k) \geq \ell(\boldsymbol{\theta}^{(k)}).$$

Es fácil notar que la dirección dada por (1), satisface

$$\mathbf{U}^\top(\boldsymbol{\theta}^{(k)})\mathbf{p}_k = \mathbf{U}^\top(\boldsymbol{\theta}^{(k)})\{-\mathbf{H}(\boldsymbol{\theta}^{(k)})\}^{-1}\mathbf{U}(\boldsymbol{\theta}^{(k)}) > 0,$$

es decir, corresponde a una dirección de ascenso.



Método de máxima verosimilitud

Sea $\mathbf{g}_k = \mathbf{U}(\boldsymbol{\theta}^{(k+1)}) - \mathbf{U}(\boldsymbol{\theta}^{(k)})$, esta clase de métodos actualiza una estimación de $\{\mathbf{H}(\boldsymbol{\theta}^{(k)})\}^{-1}$ usando un método secante, tal como el método de [Davidon-Fletcher-Powell \(DFP\)](#)

$$\mathbf{B}_{k+1}^{-1} = \mathbf{B}_k^{-1} + \frac{\mathbf{s}_k \mathbf{s}_k^\top}{\mathbf{s}_k^\top \mathbf{g}_k} - \frac{\mathbf{B}_k^{-1} \mathbf{g}_k \mathbf{g}_k^\top \mathbf{B}_k^{-1}}{\mathbf{g}_k^\top \mathbf{B}_k^{-1} \mathbf{g}_k},$$

o bien, el método de [Broyden-Fletcher-Goldfarb-Shanno \(BFGS\)](#)

$$\mathbf{B}_{k+1}^{-1} = \mathbf{B}_k^{-1} + \frac{\mathbf{s}_k \mathbf{g}_k^\top \mathbf{B}_k^{-1}}{\mathbf{s}_k^\top \mathbf{g}_k} - \frac{\mathbf{B}_k^{-1} \mathbf{g}_k \mathbf{s}_k^\top}{\mathbf{s}_k^\top \mathbf{g}_k} + \left\{ 1 + \frac{\mathbf{g}_k^\top \mathbf{B}_k^{-1} \mathbf{g}_k}{\mathbf{s}_k^\top \mathbf{g}_k} \right\} \frac{\mathbf{s}_k \mathbf{s}_k^\top}{\mathbf{s}_k^\top \mathbf{g}_k},$$

donde $\mathbf{s}_k = \boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}$.



El **método Fisher-scoring**, corresponde a un algoritmo en la clase quasi-Newton donde $-H(\theta)$ es aproximada mediante la matriz de información de Fisher.

De este modo, obtenemos el siguiente esquema iterativo:

$$\theta^{(k+1)} = \theta^{(k)} + \lambda_k \mathcal{F}^{-1}(\theta^{(k)}) U(\theta^{(k)}), \quad k = 0, 1, \dots,$$

donde $\mathcal{F}(\theta) = E\{-H(\theta)\}$.



Método de máxima verosimilitud

Cuando tenemos $\mathbf{x}_1, \dots, \mathbf{x}_n$ observaciones (vectores) independientes, tenemos que

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(\mathbf{x}_i; \boldsymbol{\theta}).$$

De este modo, $U_n(\boldsymbol{\theta}) = \sum_{i=1}^n U_i(\boldsymbol{\theta})$, con $U_i(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}_i; \boldsymbol{\theta})$. Notando que

$$\frac{1}{n} \sum_{i=1}^n U_i(\boldsymbol{\theta}) U_i^\top(\boldsymbol{\theta}) \xrightarrow{P} E\{U_n(\boldsymbol{\theta}) U_n^\top(\boldsymbol{\theta})\},$$

conforme $n \rightarrow \infty$, lleva al [Algoritmo BHHH](#) (Berndt et al., 1974),² definido como:

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + \lambda_k \left\{ \frac{1}{n} \sum_{i=1}^n U_i(\boldsymbol{\theta}^{(k)}) U_i^\top(\boldsymbol{\theta}^{(k)}) \right\}^{-1} U_n(\boldsymbol{\theta}^{(k)}), \quad k = 0, 1, \dots$$

²Berndt, Hall, Hall, Hausman (1974), *Annals of Economic and Social Measurements* 3, 653-665.



Ejemplo:

Suponga X_1, \dots, X_n variables aleatorias desde la distribución Cauchy($\theta, 1$), con densidad

$$f(x; \theta) = \frac{1}{\pi\{1 + (x - \theta)^2\}}, \quad x \in \mathbb{R}, \theta \in \mathbb{R}.$$

Así, la función de log-verosimilitud adopta la forma:

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= -n \log \pi - \log \prod_{i=1}^n \{1 + (x_i - \theta)^2\} \\ &= -n \log \pi - \sum_{i=1}^n \log(1 + (x_i - \theta)^2). \end{aligned}$$

Calculando la primera derivada, obtenemos

$$U(\theta; \mathbf{x}) = \frac{\partial \ell(\theta; \mathbf{x})}{\partial \theta} = \sum_{i=1}^n \frac{2(x_i - \theta)}{1 + (x_i - \theta)^2}.$$



Método de máxima verosimilitud

Por tanto, el estimador máximo verosímil debe satisfacer la condición de primer-orden:

$$U(\theta; \mathbf{x}) = \sum_{i=1}^n \frac{2}{1 + (x_i - \theta)^2} (x_i - \theta) = \sum_{i=1}^n \omega_i(\theta)(x_i - \theta) = 0, \quad (2)$$

donde $\omega_i(\theta) = 2/(1 + (x_i - \theta)^2)$.

Es fácil notar que la Ecuación (2) no tiene solución explícita. Para aplicar el algoritmo Newton-Raphson, calculamos

$$\frac{\partial}{\partial \theta} U(\theta; \mathbf{x}) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \omega_i(\theta)(x_i - \theta) + \sum_{i=1}^n \omega_i(\theta) \frac{\partial}{\partial \theta} (x_i - \theta),$$

como

$$\frac{\partial}{\partial \theta} \omega_i(\theta) = \frac{4}{\{1 + (x_i - \theta)^2\}^2} (x_i - \theta) = \omega_i^2(\theta)(x_i - \theta),$$

esto lleva a

$$\frac{\partial}{\partial \theta} U(\theta; \mathbf{x}) = \sum_{i=1}^n \omega_i^2(\theta)(x_i - \theta)^2 - \sum_{i=1}^n \omega_i(\theta).$$



De este modo,

$$-H(\theta; \mathbf{x}) = -\frac{\partial}{\partial \theta} U(\theta; \mathbf{x}) = \sum_{i=1}^n \omega_i(\theta) \{1 - \omega_i(\theta)(x_i - \theta)^2\}.$$

Finalmente el método Newton-Raphson, adopta la forma:

$$\begin{aligned} \theta^{(k+1)} &= \theta^{(k)} - \frac{U(\theta^{(k)}; \mathbf{x})}{H(\theta^{(k)}; \mathbf{x})} \\ &= \theta^{(k)} + \frac{\sum_{i=1}^n \omega_i(\theta^{(k)})(x_i - \theta^{(k)})}{\sum_{i=1}^n \omega_i(\theta^{(k)})\{1 - \omega_i(\theta^{(k)})(x_i - \theta^{(k)})^2\}}. \end{aligned} \quad (3)$$

Observación:

En la biblioteca [heavy](#) se encuentra una alternativa al esquema iterativo en (3) que utiliza un [Algoritmo EM](#).



Ejemplo:

Suponga X_1, \dots, X_n muestra aleatoria desde $\text{Poi}(\lambda)$. En este caso,

$$\ell(\lambda; \mathbf{x}) = \sum_{i=1}^n (x_i \log \lambda - \lambda - \log x_i!), \quad \lambda > 0,$$

despreciando aquellos términos que no dependen de λ , tenemos

$$\ell(\lambda; \mathbf{x}) = \sum_{i=1}^n (x_i \log \lambda - \lambda).$$

Considere $\phi = \log \lambda$, es decir $\lambda = e^\phi$ y note que $\phi \in \mathbb{R}$. Así,

$$\ell(\phi; \mathbf{x}) = \sum_{i=1}^n (x_i \phi - e^\phi).$$

Luego, estimamos ϕ y hacemos $\hat{\lambda}_{\text{ML}} = e^{\hat{\phi}_{\text{ML}}}$.

