

MGE-201: Suficiencia y función de verosimilitud

Felipe Osorio

felipe.osorio@uv.cl

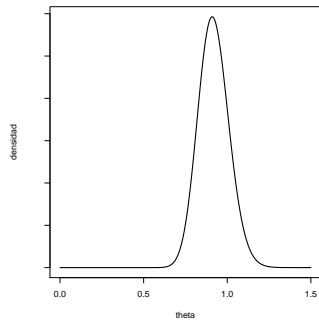
Considere X_1, \dots, X_n variables aleatorias IID desde $\text{Exp}(\theta)$, de este modo

$$\begin{aligned} f(\mathbf{x}; \theta) &= \prod_{i=1}^n \theta \exp(-\theta x_i) = \theta^n \exp\left(-\theta \sum_{i=1}^n x_i\right) \\ &= \theta^n \exp(-\theta n\bar{x}). \end{aligned}$$

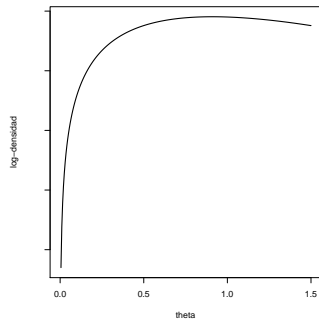
Es decir, para esta densidad conjunta **sólo** necesitamos conocer el **tamaño muestral** y la **media muestral**.

Idea:

Hemos reducido la información contenida en las n variables a una **única** estadística $T(X_1, \dots, X_n)$.



(a)



(b)

$T : \mathcal{X}^n \rightarrow \mathbb{R}$ reduce una colección de n observaciones a un único número y por tanto no puede ser inyectiva. Es decir, en general $T(X_1, \dots, X_n)$ provee **menos** información sobre θ que (X_1, \dots, X_n) .

Para algunos modelos una estadística T será igualmente informativa sobre θ que la muestra (X_1, \dots, X_n) . Tales estadísticas son llamadas **estadísticas suficientes**¹

¹Es suficiente usar T en lugar de (X_1, \dots, X_n) .

Definición 1:

Sea X_1, \dots, X_n variables aleatorias IID desde el modelo $\{P_\theta : \theta \in \Theta\}$. Una estadística $T : \mathcal{X}^n \rightarrow \mathbb{R}$ se dice **suficiente** para θ , si

$$P(X_1 \leq x_1, \dots, X_n \leq x_n | T = t),$$

no depende de θ , para todo $(x_1, \dots, x_n)^\top \in \mathbb{R}^n$ y todo $t \in \mathbb{R}$.

Ejemplo:

Suponga X_1, \dots, X_n variables aleatorias IID desde $\text{Ber}(\theta)$, donde $\theta \in (0, 1)$. Aquí $\mathcal{X} = \{0, 1\}$ mientras que $\Theta = (0, 1)$. Considere

$$T = \sum_{i=1}^n X_i,$$

sus valores son denotados como $t \in \mathcal{T} = \{0, 1, \dots, n\}$. Ahora, note que la distribución conjunta de X_1, \dots, X_n es dada por

$$p(\mathbf{x}; \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i}.$$

Por otro lado, sabemos que

$$T \sim \text{Bin}(n, \theta),$$

con probabilidad

$$p(t, \theta) = \binom{n}{t} \theta^t (1 - \theta)^{n-t}.$$

De este modo,

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n | T = t) &= \frac{P(X_1 = x_1, \dots, X_n = x_n, T = t)}{P(T = t)} \\ &= \frac{P(\{\cap_{i=1}^n X_i = x_i\} \cap \{T = t\})}{P(T = t)} = \frac{P(X_1 = x_1, \dots, X_n = x_n)}{P(T = t)} \\ &= \frac{\theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i}}{\binom{n}{t} \theta^t (1 - \theta)^{n-t}} = \frac{1}{\binom{n}{t}}. \end{aligned}$$

Es decir, conocer (X_1, \dots, X_n) además de conocer $T(X_1, \dots, X_n)$ no añade información sobre θ .

Resultado 1 (Factorización de Fisher-Neyman):

Suponga que X_1, \dots, X_n tiene densidad conjunta $f(\mathbf{x}; \boldsymbol{\theta})$, $\boldsymbol{\theta} \in \Theta$. Una estadística $T : \mathcal{X}^n \rightarrow \mathbb{R}$ es suficiente para $\boldsymbol{\theta}$ si y solo si, existe $g : \mathbb{R} \times \Theta \rightarrow \mathbb{R}$ y $h : \mathcal{X} \rightarrow \mathbb{R}$ tal que

$$f(\mathbf{x}; \boldsymbol{\theta}) = g(T(x_1, \dots, x_n); \boldsymbol{\theta})h(\mathbf{x}).$$

Demostración:

En Casella y Berger (2002, p. 276)², se presenta una demostración para el caso discreto. En el caso continuo una prueba usando Teoría de la Medida es dada en Lehmann (1986, p. 54)³

²Statistical Inference (2nd Edition). Duxbury, Pacific Grove.

³Testing Statistical Hypotheses. Wiley, New York.

Ejemplo:

Sea $\mathbf{X} = (X_1, \dots, X_n)^\top$ variables IID desde una distribución $\text{Geo}(\theta)$. De este modo, la densidad conjunta asume la forma:

$$p(\mathbf{x}; \theta) = \prod_{i=1}^n \theta(1 - \theta)^{x_i} = \theta^n (1 - \theta)^{\sum_{i=1}^n x_i},$$

para $x_i \in \{0, 1, \dots\}$. Aplicando el resultado anterior con

$$g(T(\mathbf{x}); \theta) = \theta^n (1 - \theta)^{T(\mathbf{x})}, \quad h(\mathbf{x}) = 1,$$

sigue que $T(\mathbf{x}) = \sum_{i=1}^n X_i$ es estadística suficiente.

Ejemplo:

Sea X_1, \dots, X_n una m.a.(n) desde $U(a, b)$ con $\theta = (a, b)^\top$ ($a < b$). La densidad conjunta es dada por:

$$f(\mathbf{x}; a, b) = \prod_{i=1}^n \frac{1}{b-a} I_{[a,b]}(x_i) = \frac{1}{(b-a)^n} \prod_{i=1}^n I_{[a,b]}(x_i)$$

Ahora,

$$\begin{aligned} \prod_{i=1}^n I_{[a,b]}(x_i) = 1 &\iff a \leq x_i \leq b, \forall i \\ &\iff a \leq x_{(1)} \leq x_{(n)} \leq b. \end{aligned}$$

Es decir, podemos escribir la densidad conjunta como

$$f(\mathbf{x}; a, b) = \frac{1}{(b-a)^n} I_{[a,\infty)}(x_{(1)}) I_{(-\infty,b]}(x_{(n)}).$$

De este modo, $\mathbf{T}(\mathbf{X}) = (X_{(1)}, X_{(n)})$ es suficiente para (a, b) .

Ejemplo:

Suponga $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. La densidad conjunta puede ser escrita como:

$$\begin{aligned} f(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= (2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\} \\ &= \exp\left\{\mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \frac{1}{2} \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \log |2\pi \boldsymbol{\Sigma}|\right\}, \end{aligned}$$

como $\mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} = \text{tr}(\mathbf{x} \mathbf{x}^\top \boldsymbol{\Sigma}^{-1})$, tenemos

$$f(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \exp\left\{\mathbf{T}_1^\top(\mathbf{x}) \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \frac{1}{2} \text{tr}(\mathbf{T}_2(\mathbf{x}) \boldsymbol{\Sigma}) - \frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \frac{1}{2} \log |2\pi \boldsymbol{\Sigma}|\right\},$$

con $\mathbf{T}(\mathbf{X}) = (\mathbf{T}_1(\mathbf{X}), \mathbf{T}_2(\mathbf{X}))$, es decir $\mathbf{T}_1(\mathbf{X}) = \mathbf{X}$ y $\mathbf{T}_2(\mathbf{X}) = \mathbf{X} \mathbf{X}^\top$, son estadísticas suficientes.

Ejemplo:

Sea X_1, \dots, X_n variables aleatorias IID desde $FE(\theta)$. Tenemos que,

$$f(\mathbf{x}; \theta) = \exp \left[\sum_{i=1}^n T(X_i) \eta(\theta) - nb(\theta) \right] \tilde{h}(\mathbf{x}).$$

Es decir, $\sum_{i=1}^n T(X_i)$ es estadística suficiente para θ .

Observación:

Si $T(x)$ es suficiente para θ , entonces cualquier estadístico U función 1:1 de T es suficiente para θ .

Por otro lado, sea T suficiente para θ y suponga V una función de T . Entonces, V **no** es **necesariamente suficiente** para θ .

Ejemplo:

Considere $X \sim N(\theta, 1)$ con $\theta \in \mathbb{R}$. Entonces X es suficiente para θ , mientras que $T = |X|$ no es suficiente.⁴

⁴¿Por qué?

Definición 2:

Un estadístico suficiente T , se dice **minimal**, si entre todos los estadísticos suficientes, este provee la mayor reducción de información posible.

Resultado 2:

Sea $f(\mathbf{x}; \boldsymbol{\theta})$ la densidad conjunta para la muestra $\mathbf{X} = (X_1, \dots, X_n)^\top$. Suponga que T es estadística suficiente, tal que para dos puntos cualquiera \mathbf{x} e $\mathbf{y} \in \mathcal{X}^n$ la razón

$$f(\mathbf{x}; \boldsymbol{\theta})/f(\mathbf{y}; \boldsymbol{\theta}),$$

no depende de $\boldsymbol{\theta}$ si y sólo si $T(\mathbf{x}) = T(\mathbf{y})$. Entonces $T(\mathbf{x})$ es **suficiente y minimal** para $\boldsymbol{\theta}$.

Demostración:

Considere \mathcal{K} el conjunto de todos los pares (\mathbf{x}, \mathbf{y}) para los que existe $k(\mathbf{x}, \mathbf{y}) > 0$, tal que

$$f(\mathbf{x}; \boldsymbol{\theta}) = k(\mathbf{x}, \mathbf{y})f(\mathbf{y}; \boldsymbol{\theta}), \quad \forall \boldsymbol{\theta} \in \Theta.$$

Sea S una estadística suficiente arbitraria y suponga el par (\mathbf{x}, \mathbf{y}) tal que

$$S(\mathbf{x}) = S(\mathbf{y}).$$

Por el teorema de factorización, sigue que

$$f(\mathbf{x}; \boldsymbol{\theta}) = g(S(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x}), \quad f(\mathbf{y}; \boldsymbol{\theta}) = g(S(\mathbf{y}), \boldsymbol{\theta})h(\mathbf{y}).$$

De esta manera,

$$f(\mathbf{x}; \boldsymbol{\theta}) = \frac{h(\mathbf{x})}{h(\mathbf{y})} f(\mathbf{y}; \boldsymbol{\theta}).$$

Es decir, $(\mathbf{x}, \mathbf{y}) \in \mathcal{K}$ y $T(\mathbf{x}) = T(\mathbf{y})$ lo que indica que T es una función de S .

Ejemplo:

Suponga X_1, \dots, X_n m.a.(n) desde $N(\mu, \sigma^2)$, donde μ y σ^2 son desconocidos. Sea \mathbf{x} e \mathbf{y} puntos muestrales y considere $\bar{x}, s_x^2, \bar{y}, s_y^2$ las medias y varianzas correspondientes a las muestras \mathbf{x} e \mathbf{y} . Entonces,

$$\begin{aligned} \frac{f(\mathbf{x}; \mu, \sigma^2)}{f(\mathbf{y}; \mu, \sigma^2)} &= \frac{(2\pi\sigma)^{-n/2} \exp\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\}}{(2\pi\sigma)^{-n/2} \exp\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\}} \\ &= \frac{(2\pi\sigma)^{-n/2} \exp[-\frac{1}{2\sigma^2} \{n(\bar{x} - \mu)^2 + (n-1)s_x^2\}]}{(2\pi\sigma)^{-n/2} \exp[-\frac{1}{2\sigma^2} \{n(\bar{y} - \mu)^2 + (n-1)s_y^2\}]} \\ &= \exp \left[\frac{1}{2\sigma^2} \{-n(\bar{x}^2 - \bar{y}^2) + 2n\mu(\bar{x} - \bar{y}) - (n-1)(s_x^2 - s_y^2)\} \right]. \end{aligned}$$

Esta razón no dependerá de μ y σ^2 si y sólo si $\bar{x} = \bar{y}$ y $s_x^2 = s_y^2$. De este modo, (\bar{X}, S^2) es estadística **suficiente y minimal** para (μ, σ^2) .

Ejemplo:

Considere X_1, \dots, X_n variables aleatorias IID $FE(\theta)$. Sabemos que

$$T_n(\mathbf{X}) = \sum_{i=1}^n T(X_i),$$

es estadística suficiente. Suponga muestras \mathbf{x} e \mathbf{y} , en este caso tenemos que:

$$\begin{aligned} \frac{f(\mathbf{x}; \theta)}{f(\mathbf{y}; \theta)} &= \frac{\exp\{T_n(\mathbf{x})\eta(\theta) - nb(\theta)\}\tilde{h}(\mathbf{x})}{\exp\{T_n(\mathbf{y})\eta(\theta) - nb(\theta)\}\tilde{h}(\mathbf{y})} \\ &= \frac{\tilde{h}(\mathbf{x})}{\tilde{h}(\mathbf{y})} \exp[\{T_n(\mathbf{x}) - T_n(\mathbf{y})\}\eta(\theta)], \end{aligned}$$

que es independiente de θ cuando $T_n(\mathbf{x}) = T_n(\mathbf{y})$. Por tanto, $T_n(\mathbf{X}) = \sum_{i=1}^n T(X_i)$ es **estadística suficiente minimal**.

La **información de Kullback-Leibler (KL)** entre las funciones de densidad $g(x)$ y $f(x)$ es dada por:⁵

$$I(g : f) = \int \log \left(\frac{g(x)}{f(x)} \right) g(x) \, dx = \mathbb{E}_G \left[\log \left(\frac{g(x)}{f(x)} \right) \right].$$

Propiedades de la información KL (o divergencia):

(a) $I(g : f) \geq 0$.

(b) $I(g : f) = 0 \Leftrightarrow g(x) = f(x)$ (casi en toda parte).

⁵En ocasiones anotamos $I(G : F) = \int \log(g/f) \, dG$.

Ejemplo:

Suponga que G y F están dadas, respectivamente por $N(\theta, \phi^2)$ y $N(\mu, \sigma^2)$. Entonces,

$$\begin{aligned} E_G[(X - \mu)^2] &= E_G[(X - \theta + \theta - \mu)^2] \\ &= E_G[(X - \theta)^2 + 2(X - \theta)(\theta - \mu) + (\theta - \mu)^2] \\ &= E_G[(X - \theta)^2] + (\theta - \mu)^2 = \phi^2 + (\theta - \mu)^2 \end{aligned}$$

Ahora, para

$$f(x) = (2\pi\sigma^2)^{-1/2} \exp \left\{ -\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right\},$$

sigue que

$$\begin{aligned} E_G(\log f(x)) &= E_G \left[-\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (X - \mu)^2 \right] \\ &= -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} [\phi^2 + (\theta - \mu)^2]. \end{aligned}$$

Por otro lado,

$$\mathbb{E}_G(\log g(x)) = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2}.$$

De este modo, la información KL del modelo $f(x)$ con respecto a $g(x)$ asume la forma:

$$\begin{aligned} I(g : f) &= \mathbb{E}_G(\log g(x)) - \mathbb{E}_G(\log f(x)) \\ &= \frac{1}{2} \left\{ \log \frac{\sigma^2}{\phi^2} + \frac{\phi^2 + (\theta - \mu)^2}{\sigma^2} - 1 \right\} \end{aligned}$$

Función de verosimilitud

Suponga X_1, \dots, X_n variables aleatorias IID desde una CDF desconocida $G(x)$. Asumiremos que $G(x)$ corresponde al **modelo estadístico verdadero** y sea $F(x)$ el **modelo asumido**

Supondremos también que asociadas a G y F tenemos **funciones de densidad** $g(x)$ y $f(x)$, respectivamente.

Idea:

Se desea determinar la **bondad del modelo** asumido $f(x)$ en términos de su cercanía con el modelo verdadero.

Tenemos

$$I(g : f) = \mathbb{E}_G \left[\log \frac{g(x)}{f(x)} \right] = \mathbb{E}_G[\log g(x)] - \mathbb{E}_G[\log f(x)],$$

para comparar distintos modelos competitivos basta considerar solamente el segundo término, el que es llamado **log-verosimilitud esperada**.

Observación:

Note que el cálculo de la información KL puede no ser factible pues, en general, la distribución g **no** es conocida.

Además,

$$\mathbb{E}_G[\log f(x)] = \int \log f(x) \, dG(x),$$

aún depende de la verdadera distribución. Sin embargo, podemos obtener un estimador usando en la CDF empírica \hat{G}_n basada en los datos observados X_1, \dots, X_n . Es decir,

$$\begin{aligned}\mathbb{E}_{\hat{G}_n}[\log f(x)] &= \int \log f(x) \, d\hat{G}_n(x) = \sum_{i=1}^n \hat{g}_n(x_i) \log f(x_i) \\ &= \frac{1}{n} \sum_{i=1}^n \log f(x_i).\end{aligned}$$

En efecto, de acuerdo a la Ley de los grandes números,

$$\frac{1}{n} \sum_{i=1}^n \log f(x_i) \xrightarrow{n \rightarrow \infty} \mathbb{E}_G[\log f(x)].$$

Definición 3 (Función de verosimilitud):

Para una observación \mathbf{x} fijada de un vector aleatorio \mathbf{X} con densidad $f(\cdot; \boldsymbol{\theta})$. La función de verosimilitud

$$L(\cdot; \mathbf{x}) : \Theta \rightarrow \mathbb{R}_+,$$

es definida como

$$L(\boldsymbol{\theta}; \mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}), \quad \boldsymbol{\theta} \in \Theta.$$

Observación:

La verosimilitud corresponde a la **densidad conjunta** de los datos que se desea analizar.