

Elementos de Inferencia

1.1. Suficiencia

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

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

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 $T(X_1, \dots, X_n)$.

Note que $T : \mathcal{X}^n \rightarrow \mathbb{R}$ 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* (es suficiente usar T en lugar de (X_1, \dots, X_n)).

DEFINICIÓN 1.1 (Suficiencia). 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 1.2. 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 .

TEOREMA 1.3 (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 \(2002, pág. 276\)](#), se presenta una demostración para el caso discreto. En el caso continuo una prueba usando Teoría de la Medida es dada en [Lehmann \(1986, pág. 54\)](#). \square

EJEMPLO 1.4. 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 1.5. Sea X_1, \dots, X_n una m.a.(n) desde $U(a, b)$ con $\boldsymbol{\theta} = (a, b)^\top$ ($a < b$). 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) .

1.2. Función de verosimilitud

DEFINICIÓN 1.6 (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.

EJEMPLO 1.7. Sea X_1, \dots, X_n variables aleatorias IID con distribución $N(\theta, 1)$. Entonces,

$$\begin{aligned} L(\boldsymbol{\theta}; \mathbf{x}) &= \prod_{i=1}^n (2\pi)^{-1/2} \exp \left\{ -\frac{1}{2}(x_i - \theta)^2 \right\} = (2\pi)^{-n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2 \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2 \right\}. \end{aligned}$$

Ahora,

$$\sum_{i=1}^n (x_i - \theta)^2 = \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \theta)^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \theta)^2,$$

pues $\sum_{i=1}^n (x_i - \bar{x})(\bar{x} - \theta) = 0$. De este modo,

$$\begin{aligned} L(\boldsymbol{\theta}; \mathbf{x}) &\propto \exp \left\{ -\frac{1}{2} \left[\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \theta)^2 \right] \right\} \\ &\propto \exp \left\{ -\frac{n}{2} (\bar{x} - \theta)^2 \right\}. \end{aligned}$$

Es decir, $L(\boldsymbol{\theta}; \mathbf{x})$ es proporcional a una densidad $N_1(\bar{x}, 1/n)$.

Es conveniente usar la función de log-verosimilitud dada por

$$\ell(\boldsymbol{\theta}; \mathbf{x}) = \log L(\boldsymbol{\theta}; \mathbf{x}) = \log f(\mathbf{x}; \boldsymbol{\theta}).$$

Note también que para $\boldsymbol{\theta} \in \Theta$ fijado,

$$L(\boldsymbol{\theta}; \mathbf{X}), \quad \text{y} \quad \ell(\boldsymbol{\theta}; \mathbf{X}),$$

corresponden a variables aleatorias.

EJEMPLO 1.8. Sea $X \sim N(\theta, 1)$ y $Y \sim N(2\theta, 1)$ independientes. Entonces,

$$\ell(\boldsymbol{\theta}; X) = -\frac{1}{2} \log 2\pi - \frac{1}{2}(X - \theta)^2,$$

de este modo, la variable aleatoria

$$-2\ell(\boldsymbol{\theta}; X) - \log 2\pi = (X - \theta)^2 \sim \chi^2(1).$$

Análogamente,

$$\ell(\boldsymbol{\theta}; X, Y) = \log f(X; \boldsymbol{\theta}) + \log f(Y; \boldsymbol{\theta}) = -\log 2\pi - \frac{1}{2}(X - \theta)^2 - \frac{1}{2}(Y - 2\theta)^2.$$

Es decir,

$$-2\ell(\boldsymbol{\theta}; X) - 2\log 2\pi \sim \chi^2(2).$$

OBSERVACIÓN. Considere dos conjuntos \mathbf{x}, \mathbf{y} independientes, con densidades $f_1(\mathbf{x}; \boldsymbol{\theta})$ y $f_2(\mathbf{x}; \boldsymbol{\theta})$ que comparten un parámetro común $\boldsymbol{\theta}$. Entonces la verosimilitud de los datos combinados es:

$$L(\boldsymbol{\theta}) = f_1(\mathbf{x}; \boldsymbol{\theta})f_2(\mathbf{x}; \boldsymbol{\theta}) = L_1(\boldsymbol{\theta})L_2(\boldsymbol{\theta}).$$

Además, la función de log-verosimilitud

$$\ell(\boldsymbol{\theta}) = \log L_1(\boldsymbol{\theta}) + \log L_2(\boldsymbol{\theta}) = \ell_1(\boldsymbol{\theta}) + \ell_2(\boldsymbol{\theta}).$$

El caso más simple, es para una muestra de vectores aleatorios IID. En cuyo caso, tenemos

$$L_n(\boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{x}_i; \boldsymbol{\theta}), \quad \text{y} \quad \ell_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(\mathbf{x}_i; \boldsymbol{\theta}) = \sum_{i=1}^n \ell_i(\boldsymbol{\theta}).$$

1.3. Función score e información de Fisher

A continuación definimos cantidades que surgen a partir de la log-verosimilitud. Considere los siguientes *condiciones de regularidad*

SUPUESTO A1. Las distribuciones $\{P_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \Theta\}$ tienen soporte común, de modo que el conjunto

$$A = \{\mathbf{x} : f(\mathbf{x}; \boldsymbol{\theta}) \geq 0\},$$

no depende de $\boldsymbol{\theta}$.

OBSERVACIÓN. Distribuciones pertenecientes a la FE satisfacen la condición anterior.

EJEMPLO 1.9 (Contraejemplos). Considere $X \sim \mathcal{U}(0, \theta)$, con $\theta \in (0, \infty)$ cuya densidad es

$$f(x; \theta) = \frac{1}{\theta} I_{[0, \theta]}(x).$$

También la familia de distribuciones exponencial con dos parámetros $Y \sim \text{Exp}(a, b)$,

$$f(y; a, b) = \frac{1}{b} \exp\left(-\frac{(y-a)}{b}\right) I_{[a, \infty)}(y), \quad a, b > 0.$$

SUPUESTO A2. El espacio paramétrico $\Theta \subset \mathbb{R}^p$ es un conjunto abierto.

SUPUESTO A3. Para todo $\mathbf{x} \in A$ la función de log-verosimilitud es 3-veces continuamente diferenciable con respecto a $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^\top$.

DEFINICIÓN 1.10 (Función score). Suponga las condiciones A1 a A3 para todo $\mathbf{x} \in A$, se define la *función score* como el vector de derivadas parciales de la log-verosimilitud

$$\mathbf{U}(\boldsymbol{\theta}; \mathbf{x}) = \frac{\partial \ell(\boldsymbol{\theta}; \mathbf{x})}{\partial \boldsymbol{\theta}} = \left(\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{x})}{\partial \theta_1}, \dots, \frac{\partial \ell(\boldsymbol{\theta}; \mathbf{x})}{\partial \theta_p} \right)^\top.$$

SUPUESTO A4. Suponga que existen funciones integrables $F_1(x)$, $F_2(x)$ y $H(x)$ tal que

$$\int H(x) f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} < M,$$

para $M \in \mathbb{R}$ un valor apropiado y que se satisface

$$\begin{aligned} \left| \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i} \right| &< F_1(\mathbf{x}), & \left| \frac{\partial^2 \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right| &< F_2(\mathbf{x}), \\ \left| \frac{\partial^3 \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j \partial \theta_k} \right| &< H(\mathbf{x}), & i, j, k &= 1, \dots, p. \end{aligned}$$

Esta condición implica que podemos intercambiar las operaciones de integración y diferenciación. Por ejemplo,

$$\frac{\partial}{\partial \boldsymbol{\theta}} \int_A f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \int_A \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}.$$

RESULTADO 1.11. Bajo las condiciones **A1** a **A4**, tenemos

$$\mathbb{E}_{\boldsymbol{\theta}}\{\mathbf{U}(\boldsymbol{\theta}; \mathbf{X})\} = \mathbf{0}, \quad \forall \boldsymbol{\theta} \in \Theta$$

DEMOSTRACIÓN. Tenemos

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}}\{\mathbf{U}(\boldsymbol{\theta}; \mathbf{X})\} &= \int_A \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}; \boldsymbol{\theta}) f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} \\ &= \int_A \frac{1}{f(\mathbf{x}; \boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} \\ &= \int_A \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \frac{\partial}{\partial \boldsymbol{\theta}} \int_A f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} \\ &= \mathbf{0} \end{aligned}$$

□

EJEMPLO 1.12. Considere $X \sim \mathbf{N}(\theta, 1)$. Entonces,

$$U(\theta; X) = \frac{\partial}{\partial \theta} \left[-\frac{1}{2} \log 2\pi - \frac{1}{2}(X - \theta)^2 \right] = X - \theta.$$

De este modo, $U(\theta; X) \sim \mathbf{N}(0, 1)$. Así, es directo

$$\mathbb{E}\{U(\theta; X)\} = \mathbb{E}(X - \theta) = 0.$$

DEFINICIÓN 1.13 (Matriz de información de Fisher). Suponga las condiciones **A1** a **A3**. Entonces la *matriz de información de Fisher* se define como:

$$\mathcal{F}(\boldsymbol{\theta}) = \text{Cov}_{\boldsymbol{\theta}}\{\mathbf{U}(\boldsymbol{\theta}; \mathbf{X})\} = \mathbb{E}_{\boldsymbol{\theta}}\{\mathbf{U}(\boldsymbol{\theta}; \mathbf{X})\mathbf{U}^{\top}(\boldsymbol{\theta}; \mathbf{X})\}.$$

Es decir, $\mathcal{F}(\boldsymbol{\theta})$ tiene elementos

$$\mathcal{F}_{ij}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}} \left\{ \frac{\partial}{\partial \theta_i} \ell(\boldsymbol{\theta}; \mathbf{X}) \frac{\partial}{\partial \theta_j} \ell(\boldsymbol{\theta}; \mathbf{X}) \right\}.$$

OBSERVACIÓN. En ocasiones escribimos $\mathcal{F}_X(\boldsymbol{\theta})$ pero la información **no** es aleatoria.

EJEMPLO 1.14. Sean X_1, \dots, X_n variables aleatorias $\mathbf{N}(\mu, \sigma^2)$ con σ^2 conocido. Entonces,

$$L(\mu) \propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right\},$$

así

$$\ell(\mu) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 + c,$$

con c una constante. Además,

$$U(\mu; \mathbf{X}) = \dot{\ell}(\mu) = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu).$$

De este modo,

$$\mathcal{F}_n(\mu) = \text{var}\{U(\mu; \mathbf{X})\} = \frac{1}{\sigma^4} \sum_{i=1}^n \text{var}(X_i - \mu) = \frac{n}{\sigma^2}.$$

INTERPRETACIÓN. Los datos contienen más información sobre μ si:

- (a) σ^2 es pequeño ($\sigma^2 \rightarrow 0$).
- (b) conforme n crece ($n \rightarrow \infty$).

RESULTADO 1.15. *Suponga las condiciones A1 a A4. Entonces,*

$$\mathcal{F}(\theta) = \mathbb{E}_\theta\{-\ddot{\ell}(\theta; \mathbf{X})\} = \mathbb{E}_\theta\left\{-\frac{\partial^2 \ell(\theta; \mathbf{X})}{\partial \theta \partial \theta^\top}\right\}.$$

DEMOSTRACIÓN. Tenemos que

$$\begin{aligned} \ddot{\ell}(\theta; \mathbf{x}) &= \frac{\partial \log f(\mathbf{x}; \theta)}{\partial \theta \partial \theta^\top} = \frac{\partial}{\partial \theta} \left\{ \frac{\partial}{\partial \theta^\top} \log f(\mathbf{x}; \theta) \right\} = \frac{\partial}{\partial \theta} \left\{ \frac{1}{f(\mathbf{x}; \theta)} \frac{\partial}{\partial \theta^\top} f(\mathbf{x}; \theta) \right\} \\ &= \frac{1}{f^2(\mathbf{x}; \theta)} \left\{ \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial \theta \partial \theta^\top} f(\mathbf{x}; \theta) - \frac{\partial}{\partial \theta} f(\mathbf{x}; \theta) \frac{\partial}{\partial \theta^\top} f(\mathbf{x}; \theta) \right\} \\ &= \frac{1}{f(\mathbf{x}; \theta)} \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial \theta \partial \theta^\top} - \left[\frac{1}{f(\mathbf{x}; \theta)} \frac{\partial}{\partial \theta} f(\mathbf{x}; \theta) \right] \left[\frac{1}{f(\mathbf{x}; \theta)} \frac{\partial}{\partial \theta^\top} f(\mathbf{x}; \theta) \right]. \end{aligned}$$

Por otro lado, note que

$$\begin{aligned} \mathbb{E}_\theta \left\{ \frac{1}{f(\mathbf{x}; \theta)} \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial \theta \partial \theta^\top} \right\} &= \int_A \frac{1}{f(\mathbf{x}; \theta)} \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial \theta \partial \theta^\top} f(\mathbf{x}; \theta) d\mathbf{x} \\ &= \int_A \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial \theta \partial \theta^\top} d\mathbf{x} = \frac{\partial^2}{\partial \theta \partial \theta^\top} \int_A f(\mathbf{x}; \theta) d\mathbf{x} \\ &= \mathbf{0}. \end{aligned}$$

De este modo,

$$\begin{aligned} \mathbb{E}_\theta\{-\ddot{\ell}(\theta; \mathbf{X})\} &= \mathbb{E}_\theta \left\{ \left[\frac{1}{f(\mathbf{x}; \theta)} \frac{\partial}{\partial \theta} f(\mathbf{x}; \theta) \right] \left[\frac{1}{f(\mathbf{x}; \theta)} \frac{\partial}{\partial \theta^\top} f(\mathbf{x}; \theta) \right] \right\} \\ &= \mathbb{E}_\theta \left\{ \left[\frac{\partial}{\partial \theta} \log f(\mathbf{x}; \theta) \right] \left[\frac{\partial}{\partial \theta^\top} \log f(\mathbf{x}; \theta) \right] \right\} \\ &= \mathbb{E}_\theta\{U(\theta; \mathbf{X})U^\top(\theta; \mathbf{X})\} = \text{Cov}(U(\theta; \mathbf{X})). \end{aligned}$$

□

OBSERVACIÓN. Este resultado permite obtener la matriz de información de Fisher de dos maneras equivalente en **modelos regulares** (esto es, bajo los Supuestos A1 a A4). Es decir,

$$\mathcal{F}(\theta) = \mathbb{E}_\theta\{U(\theta; \mathbf{X})U^\top(\theta; \mathbf{X})\} = \mathbb{E}_\theta \left\{ -\frac{\partial^2 \ell(\theta; \mathbf{X})}{\partial \theta \partial \theta^\top} \right\}.$$

DEFINICIÓN 1.16. La matriz

$$\mathbf{J}(\boldsymbol{\theta}; \mathbf{X}) = -\frac{\partial^2 \ell(\boldsymbol{\theta}; \mathbf{X})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top},$$

se denomina *información observada*.

OBSERVACIÓN. En efecto,

$$\mathcal{F}(\boldsymbol{\theta}) = \mathbb{E}_\theta\{\mathbf{J}(\boldsymbol{\theta}; \mathbf{X})\},$$

que también es llamada *matriz de información esperada*.

EJEMPLO 1.17. Considere $X \sim \text{Bin}(n, \theta)$ con $\theta \in (0, 1)$. De este modo,

$$\ell(\theta; x) = \log \binom{n}{x} + (n - x) \log(1 - \theta) + x \log \theta,$$

así,

$$U(\theta; x) = -\frac{n - x}{1 - \theta} + \frac{x}{\theta},$$

obteniendo la derivada de $U(\theta; x)$,

$$U'(\theta; x) = -\frac{n - x}{(1 - \theta)^2} - \frac{x}{\theta^2}.$$

Además, como $\mathbb{E}(X) = n\theta$, sigue que

$$\begin{aligned} \mathcal{F}(\theta) &= \mathbb{E}_\theta\{-U'(\theta; X)\} = \frac{n - \mathbb{E}(X)}{(1 - \theta)^2} + \frac{\mathbb{E}(X)}{\theta^2} = \frac{n - n\theta}{(1 - \theta)^2} + \frac{n\theta}{\theta} \\ &= n\left(\frac{1}{1 - \theta} + \frac{1}{\theta}\right) = \frac{n}{\theta(1 - \theta)}. \end{aligned}$$

Estimación

Suponga que tenemos X_1, \dots, X_n una muestra aleatoria desde $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$. Deseamos entender el mecanismo o modelo verdadero (basado en $\theta_0 \in \Theta$) que generó los datos. A continuación presentamos ideas básicas sobre procedimientos para seleccionar θ desde los datos (estimación) así como discutir las propiedades de tales procedimientos.

DEFINICIÓN 2.1. Una función $T : \mathcal{X}^n \rightarrow \Gamma$ es llamado un estimador (puntual). Este es usado para estimar $\gamma = g(\theta)$.

Un estimador es una regla que provee un valor plausible sobre el verdadero γ (equivalentemente θ) que generó los datos. El valor $T(\mathbf{x})$ es llamado estimación de $g(\theta)$ y corresponde a una realización de la variable aleatoria $T(\mathbf{X})$.

OBSERVACIÓN. Usualmente anotamos un estimador (y una estimación) como $\hat{\gamma} = T(X_1, \dots, X_n)$ y distinguimos el método usado como $\hat{\gamma}_{ML}$, $\hat{\gamma}_{MM}$ o $\hat{\gamma}_{LS}$.

2.1. Métodos de estimación

2.1.1. Método de los momentos. Este procedimiento **no** requiere conocer la distribución subyacente de la variable aleatoria de interés X ($\sim P_\theta \in \mathcal{P}$), sino que requiere **asumir** formas específicas para sus momentos. El objetivo es substituir estos momentos por sus contrapartes empíricas.

Para formalizar el procedimiento, considere X_1, \dots, X_n una m.a.(n) desde P_θ (unidimensional). Es decir,

$$\mathcal{P} = \{P_\theta^{\otimes n} : \theta \in \Theta\},$$

y sea

$$\mu_k = \mu_k(P_\theta) = E(X^k) = \int x^k dP_\theta = \int x^k f_X(x; \theta) dx.$$

Además, suponga que P_θ tiene momentos finitos μ_1, \dots, μ_r para algún r .

Asumiremos también que el parámetro de interés γ depende de θ a través de los momentos μ_k como:

$$\gamma = h(\mu_1(P_\theta), \dots, \mu_r(P_\theta)),$$

donde h es una función conocida. Esto lleva a la siguiente definición

DEFINICIÓN 2.2 (Estimador de momentos). Suponga X_1, \dots, X_n que sigue el modelo estadístico $\{P_\theta^{\otimes n} : \theta \in \Theta\}$. El estimador de momentos es definido como

$$\hat{\gamma}_{MM} = h(m_1, \dots, m_r),$$

donde $m_k = \hat{\mu}_k$ es el momento empírico (o muestral) de orden k , dado por

$$m_k = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

EJEMPLO 2.3. Considere X_1, \dots, X_n una muestra aleatoria desde P_θ y considere

$$\gamma = \int x f(x; \theta) dx.$$

Usando el método de momentos, tenemos que:

$$\hat{\gamma}_{MM} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

EJEMPLO 2.4. Considere que el parámetro de interés es la desviación estándar σ . Note que

$$\sigma^2 = \int (x - \mu_1)^2 f(x; \theta) dx,$$

de este modo,

$$\sigma = \sqrt{\mu_2 - \mu_1^2} = h(\mu_1, \mu_2),$$

cuyo estimador usando el método de momentos adopta la forma:

$$\hat{\sigma}_{MM} = \sqrt{m_2 - m_1^2}.$$

Note que

$$\begin{aligned} m_2 - m_1^2 &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^2 = \frac{1}{n} \left(\sum_{i=1}^n x_i^2 - \frac{1}{n} (n\bar{x})^2 \right) \\ &= \frac{1}{n} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2. \end{aligned}$$

De este modo,

$$\hat{\sigma}_{MM} = \left\{ \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right\}^{1/2},$$

es decir, $\hat{\sigma}_{MM} \neq s$ (desviación estándar muestral).

EJEMPLO 2.5. El sesgo de una variable aleatoria X con distribución F es definida como

$$\gamma = \frac{E_F(X - E_F(X))^3}{(\text{var}_F(X))^{3/2}} = \frac{\mu_3 - 3\mu_2\mu_1 + 2\mu_1^3}{(\mu_2 - \mu_1^2)^{3/2}},$$

así el estimador de momentos asume la forma

$$\hat{\gamma}_{MM} = \frac{m_3 - 3m_2m_1 + 2m_1^3}{(m_2 - m_1^2)^{3/2}} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left\{ \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right\}^{3/2}}$$

EJEMPLO 2.6. La función de distribución acumulada es definida como

$$F(t) = P_F((-\infty, t]),$$

para t fijo. Un estimador natural para la probabilidad del conjunto $(-\infty, t]$ es la frecuencia relativa,

$$\hat{F}_n(t) = \hat{F}_n(t; \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n I_{(-\infty, t]}(x_i).$$

\hat{F}_n se denomina la *función de distribución empírica*. Note que

$$m_k = \mu_k(\hat{F}_n) = \int x^k d\hat{F}_n,$$

es decir, \hat{F}_n es un estimador de momentos de F .

OBSERVACIÓN. En los ejemplos anteriores **no** hemos asumido alguna distribución específica para P_θ . Es efecto, el método de momentos se puede entender como un procedimiento libre de distribución y por tanto corresponde a un método **no paramétrico**.

EJEMPLO 2.7. Sea X_1, \dots, X_n una muestra aleatoria desde una distribución log-normal con vector de parámetros $\theta = (\mu, \sigma^2)^\top \in \mathbb{R}_+ \times \mathbb{R}_+$. Es decir, cada X_i tiene densidad

$$f(z; \mu, \sigma^2) = \frac{1}{\sigma z \sqrt{2\pi}} \exp \left\{ -\frac{(\log z - \mu)^2}{2\sigma^2} \right\}, \quad z > 0.$$

En este caso

$$\mu_1 = \exp(\mu + \sigma^2/2), \quad \mu_2 = \exp(\sigma^2) (\exp(\mu + \sigma^2/2))^2,$$

y portanto los estimadores de momentos son dados por:

$$\hat{\mu}_{\text{MM}} = 2 \log m_1 - \frac{1}{2} \log m_2, \quad \hat{\sigma}_{\text{MM}}^2 = \log m_2 - 2 \log m_1.$$

OBSERVACIÓN. Una pregunta de interés es: ¿El estimador de momentos es **único**? Considere el siguiente ejemplo.

EJEMPLO 2.8. Suponga X_1, \dots, X_n una muestra aleatoria desde $\text{Poi}(\lambda)$, $\lambda > 0$. Recuerde que

$$\mathbb{E}(X_1) = \text{var}(X_1) = \lambda.$$

De este modo, podemos considerar

$$\hat{\lambda}_{\text{MM}} = \frac{1}{n} \sum_{i=1}^n x_i \quad (= h(m_1))$$

por otro lado, otro estimador puede ser

$$\tilde{\lambda}_{\text{MM}} = h(m_1, m_2) = m_2 - m_1^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Sin embargo, ¿Cuál $\hat{\lambda}_{\text{MM}}$ o $\tilde{\lambda}_{\text{MM}}$ es mejor?¹

En la práctica, tenemos que θ es vector p -dimensional y usualmente estamos interesados en $\gamma = \theta$. De este modo tenemos que

$$\begin{aligned} \mu_1 &= g_1(\theta_1, \dots, \theta_p), \\ &\vdots \\ \mu_p &= g_p(\theta_1, \dots, \theta_p). \end{aligned}$$

Resolviendo para los p -parámetros en función de los momentos, obtenemos

$$\begin{aligned} \theta_1 &= h_1(\mu_1, \dots, \mu_p), \\ &\vdots \\ \theta_p &= h_p(\mu_1, \dots, \mu_p). \end{aligned} \tag{2.1}$$

¹Más adelante estudiaremos como comparar entre dos estimadores.

Finalmente, el estimador $\hat{\theta}_{\text{MM}}$, puede ser obtenido substituyendo en (2.1) por los momentos muestrales, es decir:

$$\begin{aligned}\hat{\theta}_1 &= h_1(m_1, \dots, m_p), \\ &\vdots \\ \hat{\theta}_p &= h_p(m_1, \dots, m_p).\end{aligned}$$

Debemos resaltar que el método de momentos requiere resolver el sistema de ecuaciones **no lineal**

$$\begin{aligned}g_1(\theta_1, \dots, \theta_p) - \mu_1 &= 0, \\ &\vdots \\ g_p(\theta_1, \dots, \theta_p) - \mu_p &= 0,\end{aligned}$$

que puede ser escrito como:

$$\Psi(\theta) = 0, \quad (2.2)$$

donde $\Psi : \mathbb{R}^p \rightarrow \mathbb{R}^p$ y por tanto $\hat{\theta}_{\text{MM}}$ corresponde a una raíz de la **ecuación de estimación** en (2.2). Note que, en general, para obtener $\hat{\theta}_{\text{MM}}$ se requiere de métodos iterativos, tal como:

$$\theta^{(s+1)} = \theta^{(s)} - \{\dot{\Psi}(\theta^{(s)})\}^{-1} \Psi(\theta^{(s)}), \quad s = 0, 1, \dots,$$

donde $\theta^{(s)}$ representa una estimación para θ en la etapa s -ésima y $\dot{\Psi}(\theta) = \partial \Psi(\theta) / \partial \theta^\top$.

OBSERVACIÓN. Una dificultad evidente del método de momentos es que $\Psi(\theta) = 0$ puede tener múltiples raíces.

2.1.2. Estimación máximo verosímil. Este es uno de los procedimientos de estimación más ampliamente usados. Es motivado por el principio de verosimilitud y los estimadores obtenidos disfrutan de buenas propiedades.

DEFINICIÓN 2.9 (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

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

RESULTADO 2.10 (Invarianza del MLE). Si $\gamma = \mathbf{g}(\theta)$ y \mathbf{g} es biyectiva. Entonces $\hat{\theta}$ es el MLE para θ si y solo si $\hat{\gamma} = \mathbf{g}(\hat{\theta})$ es el MLE para γ .

DEMOSTRACIÓN. Considere $L(\boldsymbol{\theta}; \mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta})$ y como \mathbf{g} es biyectiva tenemos que

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

Además,

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

□

Para el caso en que \mathbf{g} no sea biyectiva, considere la siguiente definición.

DEFINICIÓN 2.11. Si $\hat{\boldsymbol{\theta}}_{\text{ML}}$ es el MLE de $\boldsymbol{\theta}$ y $\boldsymbol{\gamma} = \mathbf{g}(\boldsymbol{\theta})$. Entonces el MLE de $\boldsymbol{\gamma}$ es definido como:

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

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

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

donde $\mathbf{U}(\boldsymbol{\theta}) = \partial \ell(\boldsymbol{\theta}) / \partial \boldsymbol{\theta}$ corresponde al vector score. En particular, si las ecuaciones de verosimilitud son una ecuación lineal en los parámetros, el estimador máximo verosímil puede ser expresado explícitamente.

EJEMPLO 2.12 (distribución Binomial). Sea $x \in \mathcal{X} = \{0, 1, \dots, n\}$ una realización desde $\text{Bin}(n, \theta)$. El espacio paramétrico es $\Theta = (0, 1)$ y la función de verosimilitud adopta la forma

$$L(\theta; \mathbf{x}) = \binom{n}{x} \theta^x (1 - \theta)^{n-x},$$

mientras que la log-verosimilitud es dada por

$$\ell(\theta; \mathbf{x}) = \log \binom{n}{x} + x \log \theta + (n - x) \log(1 - \theta).$$

$\hat{\boldsymbol{\theta}}_{\text{ML}}$ es solución de la ecuación

$$U(\theta; \mathbf{x}) = \dot{\ell}(\theta; \mathbf{x}) = \frac{x}{\theta} - \frac{n - x}{1 - \theta} = 0.$$

Si $x \neq 0$ y $x \neq n$ la solución existe, en cuyo caso tenemos

$$\hat{\theta}_{\text{ML}} = \frac{x}{n}.$$

EJEMPLO 2.13 (distribución Normal). 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

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

resolviendo estas ecuaciones para μ y σ^2 , sigue que

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

EJEMPLO 2.14 (distribución Uniforme). 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,²

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

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 $x_{(n)} = \max\{x_1, \dots, x_n\}$. 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 2.15 (distribución Laplace). 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

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

²Basta notar que $0 \leq x_i \leq \theta \Rightarrow x_i \leq \theta < \infty$.

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

En el siguiente ejemplo se presenta la estimación de parámetros mediante máxima verosimilitud para la clase de la familia exponencial. Por simplicidad solamente será considerado el caso de la FE 1-paramétrica.

EJEMPLO 2.16 (Familia Exponencial 1-paramétrica). Sea X_1, \dots, X_n variables aleatorias 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 la densidad conjunta es dada por

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

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 2.17 (distribución Weibull). Suponga X_1, \dots, X_n muestra aleatoria con distribución Weibull, en cuyo caso,

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

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

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

Diferenciando obtenemos las ecuaciones:

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

que corresponde a un sistema de ecuaciones no lineales y por tanto métodos iterativos son necesarios.

Debemos resaltar que los procedimientos tipo-Newton son apropiados para resolver problemas *no restringidos*. Mientras que, en general, la estimación máximo verosímil corresponde a un problema de optimización restringida. En efecto, debemos tener que $\hat{\theta}_{\text{ML}} \in \Theta$. Además, los estimadores que surgen de utilizar procedimientos de estimación restringida, suelen tener propiedades ligeramente más complejas que sus contrapartes no restringidas. El siguiente ejemplo, permite notar una forma de contornar esta dificultad mediante reparametrizar el modelo estadístico.

EJEMPLO 2.18 (distribución Poisson). 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 $\hat{\lambda}_{\text{ML}} = e^{\hat{\phi}_{\text{ML}}}$.

2.2. Propiedades de estimadores puntuales

Es frecuente contar con más de un estimador para un parámetro de interés. De este modo es requerido disponer de algún criterio que permita la comparación de diferentes estimadores. Considere las siguientes definiciones.

DEFINICIÓN 2.19 (Error Cuadrático Medio). Sea $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ un modelo estadístico para la variable X y sea T un estimador para $\gamma = g(\theta)$. El error cuadrático medio (MSE) de T es dado por

$$\text{MSE}(T, \theta) = E_\theta\{(T - g(\theta))^2\}.$$

OBSERVACIÓN. Es fácil notar que

$$\text{MSE}(T, \theta) = \{E_\theta(T) - g(\theta)\}^2 + \text{var}_\theta(T).$$

DEFINICIÓN 2.20 (Sesgo). El sesgo de un estimador T es definido como:

$$\text{bias}(T, \theta) = E_\theta(T) - g(\theta).$$

De este modo, usando la definición anterior, tenemos que:

$$\text{MSE}(T, \theta) = \{\text{bias}(T, \theta)\}^2 + \text{var}_\theta(T).$$

DEFINICIÓN 2.21 (Inssegamiento). Un estimador T para $\gamma = g(\theta)$ se dice inssegado, si

$$E_\theta(T) = g(\theta), \quad \forall \theta \in \Theta,$$

o equivalentemente,

$$\text{bias}(T, \theta) = 0, \quad \forall \theta \in \Theta.$$

Estimadores que “en promedio” están alejados de $g(\theta)$ son indeseables. Aunque en algunos casos es tolerable un sesgo pequeño. En ocasiones tenemos estimadores en que su sesgo tiende a cero conforme $n \rightarrow \infty$.

EJEMPLO 2.22. Sea X_1, \dots, X_n variables aleatorias IID con varianza finita. Suponga que $\gamma = \sigma^2$ es el parámetro de interés. Sabemos que el estimador MM es dado por:

$$\hat{\sigma}_{\text{MM}}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Además,

$$\sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n X_i^2 - n\bar{X}^2 = \mathbf{X}^\top \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \mathbf{X} = \mathbf{X}^\top \mathbf{C} \mathbf{X},$$

donde $\mathbf{C} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$. Como X_1, \dots, X_n son IID considere

$$\mathbf{E}(\mathbf{X}) = \mu \mathbf{1}_n, \quad \text{Cov}(\mathbf{X}) = \sigma^2 \mathbf{I}_n.$$

De este modo,³

$$\mathbf{E} \left\{ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right\} = \frac{1}{n} \mathbf{E}(\mathbf{X}^\top \mathbf{C} \mathbf{X}) = \frac{1}{n} \{ \sigma^2 \text{tr} \mathbf{C} + \mu^2 \mathbf{1}^\top \mathbf{C} \mathbf{1} \}.$$

Como $\text{tr} \mathbf{C} = \text{tr} \mathbf{I} - \frac{1}{n} \text{tr} \mathbf{1}\mathbf{1}^\top = n - 1$ y $\mathbf{C} \mathbf{1} = \mathbf{0}$, sigue que

$$\mathbf{E}(\hat{\sigma}_{\text{MM}}^2) = \mathbf{E} \left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right) = \left(\frac{n-1}{n} \right) \sigma^2.$$

Es decir, $\hat{\sigma}_{\text{MM}}^2$ es un estimador sesgado, y

$$\text{bias}(\hat{\sigma}_{\text{MM}}^2, \sigma^2) = \mathbf{E}(\hat{\sigma}_{\text{MM}}^2) - \sigma^2 = \frac{n-1}{n} \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n}.$$

Aunque $\lim_{n \rightarrow \infty} \text{bias}(\hat{\sigma}_{\text{MM}}^2, \sigma^2) = 0$. El “factor de corrección” $\frac{n}{n-1}$, lleva al estimador

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2,$$

que es insesgado.

Suponga adicionalmente que $X_i \sim \mathbf{N}(\mu, \sigma^2)$, para $i = 1, \dots, n$. Entonces,

$$U = \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1),$$

y sabemos que

$$\mathbf{E}(U) = n-1, \quad \text{y} \quad \text{var}(U) = 2(n-1).$$

Podemos escribir

$$U = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \implies \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{\sigma^2}{n} U,$$

luego, tenemos

$$\text{var}(\hat{\sigma}_{\text{MM}}^2) = \frac{\sigma^4}{n^2} \text{var}(U).$$

³Para \mathbf{X} vector aleatorio con $\mathbf{E}(\mathbf{X}) = \boldsymbol{\theta}$ y $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$, tenemos que $\mathbf{E}(\mathbf{X}^\top \mathbf{A} \mathbf{X}) = \text{tr} \mathbf{A} \boldsymbol{\Sigma} + \boldsymbol{\theta}^\top \mathbf{A} \boldsymbol{\theta}$.

De este modo,

$$\text{MSE}(\hat{\sigma}_{\text{MM}}^2, \sigma^2) = \left(-\frac{\sigma^2}{n}\right)^2 + \frac{\sigma^4}{n^2} \text{var}(U) = \frac{\sigma^4}{n^2} + \frac{\sigma^4}{n^2} 2(n-1) = \sigma^4 \left(\frac{2n-1}{n^2}\right).$$

Mientras que

$$U = \frac{(n-1)S^2}{\sigma^2} \implies S^2 = \frac{\sigma^2}{n-1} U.$$

Así,

$$\text{MSE}(S^2, \sigma^2) = 0 + \text{var}(S^2),$$

es decir,

$$\text{MSE}(S^2, \sigma^2) = 0 + \frac{\sigma^4}{(n-1)^2} \text{var}(U) = \frac{2(n-1)\sigma^4}{(n-1)^2} = \frac{2\sigma^4}{n-1}.$$

Finalmente,

$$\text{MSE}(S^2, \sigma^2) > \text{MSE}(\hat{\sigma}_{\text{MM}}^2, \sigma^2), \quad n > 1.$$

Es decir, aunque $\hat{\sigma}_{\text{MM}}^2$ es un estimador sesgado, este es *mejor* usando el error cuadrático medio. Note además que, bajo normalidad, $\hat{\sigma}_{\text{MM}}^2 = \hat{\sigma}_{\text{ML}}^2$.

Suponga $\Theta \subseteq \mathbb{R}^k$ y que el parámetro de interés $\boldsymbol{\gamma}$ es k -dimensional, esto es, $\boldsymbol{g} : \Theta \rightarrow \Gamma \subseteq \mathbb{R}^m$. Entonces, \boldsymbol{T} se dice un estimador insesgado, si

$$\mathbb{E}(\boldsymbol{T}) = \boldsymbol{g}(\boldsymbol{\theta}), \quad \forall \boldsymbol{\theta} \in \Theta.$$

La extensión del error cuadrático medio para el caso multiparamétrico adopta la forma

$$\begin{aligned} \text{MSE}(\boldsymbol{T}, \boldsymbol{\theta}) &= \mathbb{E}_{\boldsymbol{\theta}}\{(\boldsymbol{T} - \boldsymbol{g}(\boldsymbol{\theta}))(\boldsymbol{T} - \boldsymbol{g}(\boldsymbol{\theta}))^\top\} \\ &= \text{Cov}_{\boldsymbol{\theta}}(\boldsymbol{T}) + \{\mathbb{E}_{\boldsymbol{\theta}}(\boldsymbol{T}) - \boldsymbol{g}(\boldsymbol{\theta})\}\{\mathbb{E}_{\boldsymbol{\theta}}(\boldsymbol{T}) - \boldsymbol{g}(\boldsymbol{\theta})\}^\top. \end{aligned}$$

DEFINICIÓN 2.23. Sean \boldsymbol{T} y \boldsymbol{T}_* dos estimadores para $\boldsymbol{\gamma}$. Decimos que \boldsymbol{T}_* tiene error cuadrático medio más pequeño que \boldsymbol{T} si

$$\boldsymbol{u}^\top (\text{MSE}(\boldsymbol{T}_*, \boldsymbol{\theta}) - \text{MSE}(\boldsymbol{T}, \boldsymbol{\theta})) \boldsymbol{u} \leq 0, \quad \forall \boldsymbol{u} \in \mathbb{R}^m,$$

y escribimos

$$\text{MSE}(\boldsymbol{T}_*, \boldsymbol{\theta}) \leq \text{MSE}(\boldsymbol{T}, \boldsymbol{\theta}).$$

En general, evaluar la condición dada por la definición anterior puede ser difícil. Esto ha motivado la introducción de algunos criterios más simples para comparar entre diferentes estimadores. En efecto, decimos que \boldsymbol{T}_* es *T-óptimo*, si

$$\text{tr MSE}(\boldsymbol{T}_*, \boldsymbol{\theta}) \leq \text{tr MSE}(\boldsymbol{T}, \boldsymbol{\theta}), \quad (2.3)$$

mientras que \boldsymbol{T}_* se dice *D-óptimo*, si satisface

$$\det \text{MSE}(\boldsymbol{T}_*, \boldsymbol{\theta}) \leq \det \text{MSE}(\boldsymbol{T}, \boldsymbol{\theta}). \quad (2.4)$$

El criterio dado en la Definición 2.23 también es conocido como *M-optimalidad*. Debemos destacar que en ocasiones el error cuadrático medio es definido como

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}}[\|\boldsymbol{T} - \boldsymbol{g}(\boldsymbol{\theta})\|^2] &= \|\mathbb{E}_{\boldsymbol{\theta}}(\boldsymbol{T}) - \boldsymbol{g}(\boldsymbol{\theta})\|^2 + \sum_{j=1}^m \text{var}(T_j) \\ &= \|\text{bias}(\boldsymbol{T}, \boldsymbol{\theta})\|^2 + \text{tr Cov}_{\boldsymbol{\theta}}(\boldsymbol{T}), \end{aligned}$$

que corresponde al criterio de *T-optimalidad*. Lamentablemente, es muy poco frecuente encontrar un estimador que *siempre* (es decir, para todo $\boldsymbol{\theta} \in \Theta$) sea mejor.