

# MAT-206: Sesión 12, Funciones de inferencia

**Felipe Osorio**

[fosorios.mat.utfsm.cl](mailto:fosorios.mat.utfsm.cl)

Departamento de Matemática, UTFSM



Sea  $\mathbf{Y}$  vector aleatorio  $N$ -dimensional desde un modelo estadístico

$$\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}, \quad \Theta \in \mathbb{R}^k,$$

y  $\mathcal{Y}$  es el espacio muestral.

## Definición 1:

Una función  $\Psi : \Theta \times \mathcal{Y} \rightarrow \mathbb{R}^k$  tal que  $\Psi(\theta; \cdot)$  es medible para todo  $\theta \in \Theta$  se dice una **función de inferencia**.

## Observación:

Para una función de inferencia  $\Psi$  y una muestra  $\mathbf{Y} \in \mathcal{Y}$  dadas, es posible obtener un estimador  $\hat{\theta} = \hat{\theta}(\mathbf{Y})$  como **solución de la ecuación**

$$\Psi(\theta; \mathbf{Y}) = \mathbf{0}.$$



## *Ejemplo (Función score):*

Denote por  $\ell(\theta) = \log p(\theta)$  a la función de log-verosimilitud para  $\theta$ . La **función score**

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

corresponde a una función de inferencia.

Note que, aquél valor  $\hat{\theta}_{\text{ML}}$  que maximiza  $\ell(\theta)$  y es una solución de la ecuación

$$U(\theta) = 0,$$

sobre  $\Theta$ , es llamado el **estimador máximo verosímil** para  $\theta$ .



## Definición 2 (Función de inferencia regular):

Una función de inferencia  $\Psi$  se dice **regular** para todo  $\theta \in \Theta$  si satisface las siguientes condiciones:

- (i) Tiene segundo momento finito.
- (ii) Es **insesgada**, esto es,  $E_{\theta}\{\Psi(\theta)\} = 0$ .
- (iii) La **matriz de sensibilidad**

$$S_{\Psi}(\theta) = E_{\theta} \left\{ - \frac{\partial \Psi(\theta)}{\partial \theta^{\top}} \right\},$$

es no singular para todo  $\theta$ .



## Definición 3 (Matriz de variabilidad):

Considere  $\Psi$  función de inferencia regular la **matriz de variabilidad**, es definida como:

$$V_{\Psi}(\theta) = \text{Cov}_{\theta}(\Psi(\theta)) = E_{\theta}\{\Psi(\theta)\Psi^{\top}(\theta)\},$$

## Definición 4 (Información de Godambe):

Para  $\Psi$  función de inferencia regular, la **matriz de información de Godambe** es dada por:

$$G_{\Psi}(\theta) = S_{\Psi}^{\top}(\theta)V_{\Psi}^{-1}(\theta)S_{\Psi}(\theta).$$



## *Observación:*

Note que, para  $C(\theta)$  matriz no estocástica

$$\Phi(\theta; \mathbf{Y}) = C(\theta)\Psi(\theta; \mathbf{Y}),$$

también es una función de inferencia. En cuyo caso anotamos  $\Phi \sim \Psi$ .

## *Observación:*

Asuma las condiciones A1 a A4. De este modo,  $E\{U(\theta)\} = \mathbf{0}$  y tenemos que

$$\mathcal{F}(\theta) = \text{Cov}\{U(\theta)\},$$

corresponde a la matriz de **información de Fisher**.



## Definición 5 (Función de estimación óptima):

Sea  $\Psi$  y  $\Phi$  funciones de inferencia regulares para el vector de parámetros  $\theta$ . Si se satisface que

$$G_{\Psi}^{-1}(\theta) - G_{\Phi}^{-1}(\theta) \geq 0,$$

para cualquier  $\Psi$ , entonces se dice que  $\Phi$  es **función de inferencia óptima**.

## Resultado 1 (Desigualdad de Godambe):

Asuma que  $\Psi$  es función de inferencia regular. Entonces

$$G_{\Psi}^{-1}(\theta) - \mathcal{F}^{-1}(\theta) \geq 0,$$

para todo  $\theta \in \Theta$ , donde la igualdad se cumple sólo si  $\Psi \sim U$ , la función score.



## Esbozo de la demostración:

Considere el caso unidimensional ( $k = 1$ ) y sea  $U_n(\theta) = d \log p(\mathbf{y}; \theta) / d\theta$ . Como  $\Psi(\theta)$  es regular, tenemos

$$\int \Psi(\theta) p(\mathbf{y}; \theta) d\mathbf{y} = 0,$$

de ahí que

$$\begin{aligned} 0 &= \frac{d}{d\theta} \int \Psi(\theta) p(\mathbf{y}; \theta) d\mathbf{y} = \int \frac{d}{d\theta} \{ \Psi(\theta) p(\mathbf{y}; \theta) \} d\mathbf{y} \\ &= \int \frac{d}{d\theta} \Psi(\theta) p(\mathbf{y}; \theta) d\mathbf{y} + \int \Psi(\theta) \frac{d}{d\theta} p(\mathbf{y}; \theta) d\mathbf{y} \end{aligned}$$

Por otro lado,

$$\frac{d}{d\theta} p(\mathbf{y}; \theta) = \frac{p(\mathbf{y}; \theta)}{p(\mathbf{y}; \theta)} \frac{d}{d\theta} p(\mathbf{y}; \theta) = p(\mathbf{y}; \theta) \frac{d}{d\theta} \log p(\mathbf{y}; \theta) = p(\mathbf{y}; \theta) U_n(\theta).$$





De este modo,

$$0 = \int \frac{d}{d\theta} \Psi(\theta) p(\mathbf{y}; \theta) d\mathbf{y} + \int \Psi(\theta) U_n(\theta) p(\mathbf{y}; \theta) d\mathbf{y}.$$

Notando que

$$\int \Psi(\theta) U_n(\theta) p(\mathbf{y}; \theta) d\mathbf{y} = \text{Cov}(\Psi(\theta), U_n(\theta)).$$

Es decir, tenemos que

$$E_{\theta}^2 \left( - \frac{d \Psi(\theta)}{d \theta} \right) = \text{Cov}^2 (\Psi(\theta), U_n(\theta)) \leq E_{\theta}(\Psi^2(\theta)) E_{\theta}(U_n^2(\theta)).$$

Por lo tanto,

$$E_{\theta}(U_n^2(\theta)) \geq \frac{E_{\theta}^2(-d \Psi(\theta) / d \theta)}{E_{\theta}(\Psi^2(\theta))}.$$

Finalmente,

$$\mathcal{F}^{-1}(\theta) \leq G_{\Psi}^{-1}(\theta).$$



## Definición 6 (Algoritmo Newton-scoring):

Considere la expansión en series de Taylor de  $\Psi(\theta)$  en torno de  $\theta_0$ , tenemos

$$\Psi(\theta) \approx \Psi(\theta_0) + \frac{\partial \Psi(\theta)}{\partial \theta^\top} \Big|_{\theta=\theta_0} (\theta - \theta_0),$$

como  $\Psi(\hat{\theta}) = 0$  y substituyendo  $\dot{\Psi}(\theta_0)$  por  $S_\Psi(\theta_0)$ , sigue que

$$\hat{\theta} = \theta_0 + S_\Psi^{-1}(\theta_0) \Psi(\theta_0),$$

esto sugiere considerar:

$$\theta^{(t+1)} = \theta^{(t)} + S_\Psi^{-1}(\theta^{(t)}) \Psi(\theta^{(t)}),$$

para llevar a cabo la estimación de parámetros.

## Observación:

Este procedimiento fue propuesto por Jørgensen y Knudsen (2004)<sup>1</sup> quienes lo denominaron **algoritmo Newton-scoring**.

---

<sup>1</sup>Scandinavian Journal of Statistics **31**, 93-114.



## Resultado 2:

Sea  $\{\hat{\theta}_n\}_{n \geq 1}$  una secuencia de raíces de las ecuaciones de estimación

$$\Psi_n(\theta) = \sum_{i=1}^n \Psi_i(\theta; Y_i) = 0,$$

desde la condición de inesgamiento se tiene la **consistencia**  $\hat{\theta}_n \xrightarrow{P} \theta_0$ .

## Definición 7 (Optimalidad de Crowder):

Considere la función de inferencia regular (aditiva):

$$\Psi_n(\theta) = \sum_{i=1}^n C_i(\theta) \Psi_i(\theta; Y_i), \quad \theta \in \Theta,$$

donde  $C_i(\theta)$  es una matriz no aleatoria de  $\theta$ , tal que la secuencia  $\{\hat{\theta}_n\}_{n \geq 1}$  es consistente. Entonces la **función de inferencia optimal** está definida por

$$C_i(\theta) = E_{\theta}^{\top} \{-\dot{\Psi}_i(\theta)\} \text{Cov}_{\theta}^{-1} \{\Psi_i(\theta)\}.$$



## Supuestos

**C1:**  $\Psi_n(\theta) \rightarrow \mathbf{0}$  con probabilidad 1.

**C2:** Existe una vecindad de  $\theta_0$  tal que, con prob. 1,  $\Psi_n(\theta)$  es diferenciable y  $\dot{\Psi}_n(\theta)$  converge a un límite no estocástico que es no singular en  $\theta_0$ .

**C3:**  $\frac{1}{\sqrt{n}} \Psi_n(\theta) \xrightarrow{D} N_k(\mathbf{0}, V(\theta_0))$ .

## Resultado 3 (Normalidad asintótica):

Si  $\hat{\theta}_n \xrightarrow{P} \theta_0$  y  $\Psi_n(\hat{\theta}_n) = \mathbf{0}$  con probabilidad 1. Entonces, bajo los supuestos C1 a C3, tenemos que

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N_p(\mathbf{0}, G_{\Psi}^{-1}(\theta_0)).$$

## Demostración:

Detalles en Yuan y Jennrich (1998)<sup>2</sup>

---

<sup>2</sup>Journal of Multivariate Analysis **65**, 245-260.



### *Ejemplo (Ecuaciones de estimación generalizadas, GEE):*

Liang y Zeger (1986)<sup>3</sup> propusieron GEE como un método para análisis de **datos con estructura longitudinal**.

Considere  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{im_i})^\top$  vector de respuestas  $m_i$ -dimensional asociada al  $i$ -ésimo individuo. Es asumido que

$$E(Y_{ij}) = \mu_{ij}(\boldsymbol{\beta}), \quad \text{var}(Y_{ij}) = \phi^{-1}V(\mu_{ij}),$$

donde  $V(\mu)$  es función de varianza y  $\phi > 0$ . La estimación de  $\boldsymbol{\beta} \in \mathbb{R}^p$  se lleva a cabo mediante resolver el **sistema de ecuaciones**:

$$\boldsymbol{\Psi}(\boldsymbol{\beta}) := \sum_{i=1}^n \left( \frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}^\top} \right)^\top \{ \mathbf{A}_i^{1/2} \mathbf{R}_i(\boldsymbol{\alpha}) \mathbf{A}_i^{1/2} \}^{-1} (\mathbf{Y}_i - \boldsymbol{\mu}_i) = \mathbf{0},$$

con  $\mathbf{A}_i = \text{diag}(V(\mu_{i1}), \dots, V(\mu_{im_i}))$  y  $\mathbf{R}_i(\boldsymbol{\alpha})$  es conocida como matriz de correlación de trabajo.

---

<sup>3</sup>Biometrika **73**, 13-22.

## *Ejemplo (Quasi-verosimilitud):*

Suponga que  $E(Y) = \mu$  y  $\text{var}(Y) = \phi^{-1}V(\mu)$ . Wedderburn (1974)<sup>4</sup> mediante una analogía con la función score para la familia exponencial definió

$$Q(\mu; Y) = \int_y^\mu \phi \frac{y - t}{V(\mu)} dt.$$

Para **observaciones dependientes** McCullagh y Nelder (1989)<sup>5</sup> consideran

$$Q_i(\mu_i; Y_i) = -\frac{1}{\phi} (Y_i - \mu_i)^\top \left( \int_0^1 s \{V_i(t(s))\}^{-1} ds \right) (Y_i - \mu_i),$$

donde la integral está definida a lo largo de la recta  $t(s) = Y_i + (Y_i - \mu_i)s$ , para  $0 \leq s \leq 1$ . Una **función de inferencia** para  $\beta$  se obtiene como la primera derivada de

$$Q(\mu; Y) = \sum_{i=1}^n Q_i(\mu_i; Y_i).$$

---

<sup>4</sup>Biometrika **61**, 439-447.

<sup>5</sup>Generalized Linear Models, Chapman & Hall, London.

### *Ejemplo (Función de inferencia cuadrática, QIF):*

Suponga que  $\mathbf{R}^{-1} = \sum_{k=1}^m a_k \mathbf{M}_k$ , con  $\mathbf{M}_1, \dots, \mathbf{M}_m$ , matrices conocidas, mientras que  $a_1, \dots, a_m$  son constantes desconocidas. Haciendo  $\Phi_n(\beta) = \frac{1}{n} \sum_{i=1}^n \Phi_i(\beta)$  con

$$\Phi_i(\beta) = \begin{pmatrix} \mathbf{D}_i^\top \mathbf{A}_i^{-1/2} \mathbf{M}_1 \mathbf{A}_i^{-1/2} (\mathbf{Y}_i - \boldsymbol{\mu}_i) \\ \vdots \\ \mathbf{D}_i^\top \mathbf{A}_i^{-1/2} \mathbf{M}_m \mathbf{A}_i^{-1/2} (\mathbf{Y}_i - \boldsymbol{\mu}_i) \end{pmatrix},$$

donde  $\mathbf{D}_i = \partial \boldsymbol{\mu}_i / \partial \beta^\top$  para  $i = 1, \dots, n$ .

Podemos estimar  $\beta$ , en el contexto de GEE, mediante minimizar la QIF:

$$Q_n(\beta) = n \Phi_n^\top(\beta) \mathbf{C}_n^{-1}(\beta) \Phi_n(\beta),$$

con  $\mathbf{C}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \Phi_i(\beta) \Phi_i^\top(\beta)$ .

