

NOTAS DE CLASE :

Elementos de Inferencia Estadística

Felipe Osorio

DEPARTAMENTO DE MATEMÁTICA, UNIVERSIDAD TÉCNICA FEDERICO SANTA MARÍA

Índice general

Capítulo 1. Preliminares	1
1.1. Vectores Aleatorios	1
1.2. Operadores de esperanza y covarianza	3
1.3. Independencia de vectores aleatorios	8
1.4. Cambios de variable	9
1.5. Modelo estadístico	9
1.6. Familia exponencial	12
1.7. Familia de posición-escala	17
1.8. Familia de distribuciones de contornos elípticos	18
Capítulo 2. Elementos de Inferencia	25
2.1. Suficiencia	25
2.2. Función de verosimilitud	27
2.3. Función score e información de Fisher	30
Capítulo 3. Estimación	35
3.1. Métodos de estimación	35
3.2. Propiedades de estimadores puntuales	45
3.3. Propiedades Asintóticas	50
Capítulo 4. Intervalos y Regiones de Confianza	55
4.1. Método de la Cantidad Pivotal	56
4.2. Intervalos de Confianza Asintóticos	57
4.3. Regiones de Confianza Asintóticas	58
Bibliografía	61

Preliminares

1.1. Vectores Aleatorios

El propósito de esta sección es introducir algunas propiedades elementales de vectores aleatorios útiles a lo largo de este curso. Se asume que el lector es familiar con el concepto de variable aleatoria unidimensional.

Un vector aleatorio n -dimensional \mathbf{X} es una función (medible) desde el espacio de probabilidad Ω a \mathbb{R}^n , esto es

$$\mathbf{X} : \Omega \rightarrow \mathbb{R}^n.$$

Por convención asumiremos que el vector aleatorio $\mathbf{X} = (X_1, \dots, X_n)^\top$ es un vector columna.

DEFINICIÓN 1.1 (Función de distribución). Para \mathbf{X} distribuido en \mathbb{R}^n , la *función de distribución* de \mathbf{X} es una función $F : \mathbb{R}^n \rightarrow [0, 1]$, tal que

$$F(\mathbf{x}) = P(\mathbf{X} \leq \mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^n \quad (1.1)$$

y denotamos $\mathbf{X} \sim F$ o $\mathbf{X} \sim F_X$.

La función en (1.1) debe ser entendida como

$$F(\mathbf{x}) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n),$$

que corresponde a la probabilidad del evento $\bigcap_{k=1}^n \{X_k \leq x_k\}$.

PROPIEDAD 1.2. La función de distribución acumulada tiene las siguientes propiedades:

- (a) $F(\mathbf{x})$ es función monótona creciente y continua a la derecha en cada uno de los componentes de \mathbf{X} ,
- (b) $0 \leq F(\mathbf{x}) \leq 1$,
- (c) $F(-\infty, x_2, \dots, x_n) = \dots = F(x_1, \dots, x_{n-1}, -\infty) = 0$,
- (d) $F(+\infty, \dots, +\infty) = 1$.

Sea F la función de distribución del vector aleatorio \mathbf{X} . Entonces, existe una función no-negativa f tal que

$$F(\mathbf{x}) = \int_{-\infty}^{\mathbf{x}} f(\mathbf{u}) \, d\mathbf{u}, \quad \mathbf{x} \in \mathbb{R}^n,$$

en este caso decimos que \mathbf{X} es un vector aleatorio continuo con *función de densidad* f . Por el teorema fundamental del Cálculo, tenemos que

$$f(\mathbf{x}) = \frac{\partial^n F(\mathbf{x})}{\partial x_1 \cdots \partial x_n}.$$

Además, considere $\bar{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$, para \mathbf{x}, \mathbf{y} vectores en $\bar{\mathbb{R}}^n$, entonces

$$\mathbf{x} \leq \mathbf{y} \quad \text{esto es,} \quad x_i \leq y_i, \quad \text{para } i = 1, \dots, n.$$

Esto permite definir un rectángulo n -dimensional en \mathbb{R}^n como

$$I = (\mathbf{a}, \mathbf{b}] = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a} < \mathbf{x} \leq \mathbf{b}\}$$

para todo $\mathbf{a}, \mathbf{b} \in \bar{\mathbb{R}}^n$. Entonces, también por el teorema fundamental del Cálculo, tenemos que si

$$f(\mathbf{x}) = \frac{\partial^n F(\mathbf{x})}{\partial x_1 \cdots \partial x_n}.$$

existe y es continua (casi en toda parte) sobre un rectángulo I , entonces

$$P(\mathbf{x} \in A) = \int_A f(\mathbf{x}) d\mathbf{x}, \quad \forall A \subset I.$$

Naturalmente la función de densidad debe satisfacer

$$\int_{\mathbb{R}^n} f(\mathbf{x}) d\mathbf{x} = 1.$$

Considere el vector aleatorio n -dimensional \mathbf{X} particionado como $\mathbf{X} = (\mathbf{X}_1^\top, \mathbf{X}_2^\top)^\top$ donde \mathbf{X}_1 y \mathbf{X}_2 son vectores $n_1 \times 1$ y $n_2 \times 1$, respectivamente, con $n = n_1 + n_2$. Tenemos que $\mathbf{X}_i \sim F_i$, $i = 1, 2$, de este modo \mathbf{X} se denomina la *conjunta* de $\mathbf{X}_1, \mathbf{X}_2$ mientras que los \mathbf{X}_1 y \mathbf{X}_2 son llamados *marginales* de \mathbf{X} .

Note que, las funciones de distribución marginal pueden ser recuperadas desde la distribución conjunta mediante

$$F_1(\mathbf{s}) = F(\mathbf{s}, +\infty), \quad F_2(\mathbf{t}) = F(+\infty, \mathbf{t}), \quad \forall \mathbf{s} \in \mathbb{R}^{n_1}, \mathbf{t} \in \mathbb{R}^{n_2}.$$

Cuando \mathbf{X} es absolutamente continua con función de densidad $f(\mathbf{x}) = f(\mathbf{x}_1, \mathbf{x}_2)$, entonces la función de densidad de \mathbf{X}_i también es absolutamente continua y puede ser obtenida como

$$f_1(\mathbf{s}) = \int_{\mathbb{R}^{n_2}} f(\mathbf{s}, \mathbf{u}) d\mathbf{u}, \quad f_2(\mathbf{t}) = \int_{\mathbb{R}^{n_1}} f(\mathbf{u}, \mathbf{t}) d\mathbf{u}, \quad \forall \mathbf{s} \in \mathbb{R}^{n_1}, \mathbf{t} \in \mathbb{R}^{n_2},$$

el resultado anterior es análogo para el caso de distribuciones discretas. Si \mathbf{X} es absolutamente continuo y $f_1(\mathbf{x}_1) > 0$, entonces la *densidad condicional* de \mathbf{X}_2 dado $\mathbf{X}_1 = \mathbf{x}_1$ es

$$f_{X_2|X_1=\mathbf{x}_1}(\mathbf{u}) = \frac{f_X(\mathbf{x}_1, \mathbf{u})}{f_1(\mathbf{x}_1)},$$

con función de distribución de \mathbf{X}_2 condicional a $\mathbf{X}_1 = \mathbf{x}_1$ dada por

$$F_{X_2|X_1=\mathbf{x}_1}(\mathbf{u}) = \int_{-\infty}^{\mathbf{u}} f_{X_2|X_1=\mathbf{x}_1}(\mathbf{t}) d\mathbf{t},$$

tenemos además que

$$f_{X_2|X_1=\mathbf{x}_1}(\mathbf{u}) = \frac{f_X(\mathbf{x}_1, \mathbf{u})}{\int_{\mathbb{R}^{n_2}} f_X(\mathbf{x}_1, \mathbf{t}) d\mathbf{t}}.$$

1.2. Operadores de esperanza y covarianza

Considere $\mathbf{X} = (X_1, \dots, X_n)^\top$ vector aleatorio n -dimensional con función de densidad f . Entonces la esperanza de cualquier función \mathbf{g} de \mathbf{X} está dada por

$$\mathbb{E}(\mathbf{g}(\mathbf{X})) = \int_{\mathbb{R}^n} \mathbf{g}(\mathbf{t}) f(\mathbf{t}) d\mathbf{t},$$

siempre que la integral (n -dimensional) exista.

Más generalmente, sea $\mathbf{Z} = (Z_{ij})$ una función matricial $m \times n$, entonces podemos definir el operador de esperanza de una matriz aleatoria como

$$\mathbb{E}(\mathbf{Z}(\mathbf{X})) = (\mathbb{E}(Z_{ij})), \quad Z_{ij} = Z_{ij}(\mathbf{X}). \quad (1.2)$$

De la definición en (1.2) se desprenden una serie de resultados útiles con relación al operador de esperanza. Por ejemplo, sea $\mathbf{A} = (a_{ij})$ una matriz de constantes, entonces

$$\mathbb{E}(\mathbf{A}) = \mathbf{A}.$$

RESULTADO 1.3. Sea $\mathbf{A} = (a_{ij})$, $\mathbf{B} = (b_{ij})$ y $\mathbf{C} = (c_{ij})$ matrices de constantes $l \times m$, $n \times p$ y $l \times p$, respectivamente. Entonces

$$\mathbb{E}(\mathbf{AZB} + \mathbf{C}) = \mathbf{A} \mathbb{E}(\mathbf{Z}) \mathbf{B} + \mathbf{C}.$$

DEMOSTRACIÓN. Sea $\mathbf{Y} = \mathbf{AZB} + \mathbf{C}$, entonces

$$Y_{ij} = \sum_{r=1}^m \sum_{s=1}^n a_{ir} Z_{rs} b_{sj} + c_{ij},$$

de este modo

$$\begin{aligned} \mathbb{E}(\mathbf{AZB} + \mathbf{C}) &= (\mathbb{E}(Y_{ij})) = \left(\sum_{r=1}^m \sum_{s=1}^n a_{ir} \mathbb{E}(Z_{rs}) b_{sj} + c_{ij} \right) \\ &= \mathbf{A} \mathbb{E}(\mathbf{Z}) \mathbf{B} + \mathbf{C}. \end{aligned}$$

□

Un caso particular importante corresponde a la esperanza de una transformación lineal. Considere el vector aleatorio n -dimensional, $\mathbf{Y} = \mathbf{AX}$, donde \mathbf{X} es vector aleatorio $m \times 1$, entonces $\mathbb{E}(\mathbf{AX}) = \mathbf{A} \mathbb{E}(\mathbf{X})$. Esta propiedad puede ser extendida para sumas de vectores aleatorios, como

$$\mathbb{E} \left(\sum_i \mathbf{A}_i \mathbf{X}_i \right) = \sum_i \mathbf{A}_i \mathbb{E}(\mathbf{X}_i),$$

de manera similar tenemos que

$$\mathbb{E} \left(\sum_i \alpha_i \mathbf{Z}_i \right) = \sum_i \alpha_i \mathbb{E}(\mathbf{Z}_i),$$

donde α_i son constantes y los \mathbf{Z}_i son matrices aleatorias.

DEFINICIÓN 1.4 (Matriz de covarianza). Sean \mathbf{X} e \mathbf{Y} vectores aleatorios m y n -dimensionales, respectivamente. Se define la *matriz de covarianza* entre \mathbf{X} e \mathbf{Y} como la matriz $m \times n$,

$$\text{Cov}(\mathbf{X}, \mathbf{Y}) = (\text{Cov}(X_i, Y_j)).$$

Podemos apreciar, a partir de la definición de covarianza que

$$\text{Cov}(\mathbf{X}, \mathbf{Y}) = \mathbb{E}\{(\mathbf{X} - \mathbb{E}(\mathbf{X}))(\mathbf{Y} - \mathbb{E}(\mathbf{Y}))^\top\}.$$

En efecto, sean $\boldsymbol{\mu} = \mathbb{E}(\mathbf{X})$ y $\boldsymbol{\eta} = \mathbb{E}(\mathbf{Y})$. Entonces,

$$\begin{aligned}\text{Cov}(\mathbf{X}, \mathbf{Y}) &= (\text{Cov}(X_i, Y_j)) = (\mathbb{E}(X_i - \mu_i)(Y_j - \eta_j)) \\ &= \mathbb{E}([(X_i - \mu_i)(Y_j - \eta_j)]) = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\eta})^\top].\end{aligned}$$

Tenemos además el siguiente resultado

$$\begin{aligned}\text{Cov}(\mathbf{X}, \mathbf{Y}) &= \mathbb{E}\{(\mathbf{X} - \mathbb{E}(\mathbf{X}))(\mathbf{Y} - \mathbb{E}(\mathbf{Y}))^\top\} \\ &= \mathbb{E}(\mathbf{X}\mathbf{Y}^\top - \mathbb{E}(\mathbf{X})\mathbf{Y}^\top - \mathbf{X}\mathbb{E}^\top(\mathbf{Y}) + \mathbb{E}(\mathbf{X})\mathbb{E}^\top(\mathbf{Y})) \\ &= \mathbb{E}(\mathbf{X}\mathbf{Y}^\top) - \mathbb{E}(\mathbf{X})\mathbb{E}^\top(\mathbf{Y}).\end{aligned}$$

Se define la *matriz de dispersión (varianza)*, como $\text{Cov}(\mathbf{X}) = \text{Cov}(\mathbf{X}, \mathbf{X})$. De este modo, tenemos

$$\text{Cov}(\mathbf{X}) = (\text{Cov}(X_i, X_j)) = \mathbb{E}\{(\mathbf{X} - \mathbb{E}(\mathbf{X}))(\mathbf{X} - \mathbb{E}(\mathbf{X}))^\top\},$$

y, de la misma manera que para el caso de la matriz de covarianza,

$$\text{Cov}(\mathbf{X}) = \mathbb{E}(\mathbf{X}\mathbf{X}^\top) - \mathbb{E}(\mathbf{X})\mathbb{E}^\top(\mathbf{X}).$$

EJEMPLO 1.5. Sea \mathbf{a} vector de constantes $n \times 1$, entonces

$$\text{Cov}(\mathbf{X} - \mathbf{a}) = \text{Cov}(\mathbf{X}).$$

En efecto, note que

$$\mathbf{X} - \mathbf{a} - \mathbb{E}(\mathbf{X} - \mathbf{a}) = \mathbf{X} - \mathbb{E}(\mathbf{X}),$$

por tanto, tenemos

$$\text{Cov}(\mathbf{X} - \mathbf{a}, \mathbf{X} - \mathbf{a}) = \text{Cov}(\mathbf{X}, \mathbf{X})$$

RESULTADO 1.6. Si \mathbf{X} e \mathbf{Y} son vectores aleatorios m y n -dimensionales, respectivamente y $\mathbf{A} \in \mathbb{R}^{l \times m}$, $\mathbf{B} \in \mathbb{R}^{p \times n}$, entonces

$$\text{Cov}(\mathbf{A}\mathbf{X}, \mathbf{B}\mathbf{Y}) = \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{Y}) \mathbf{B}^\top.$$

DEMOSTRACIÓN. Sean $\mathbf{U} = \mathbf{A}\mathbf{X}$ y $\mathbf{V} = \mathbf{B}\mathbf{Y}$, entonces

$$\begin{aligned}\text{Cov}(\mathbf{A}\mathbf{X}, \mathbf{B}\mathbf{Y}) &= \text{Cov}(\mathbf{U}, \mathbf{V}) = \mathbb{E}\{(\mathbf{U} - \mathbb{E}(\mathbf{U}))(\mathbf{V} - \mathbb{E}(\mathbf{V}))^\top\} \\ &= \mathbb{E}\{(\mathbf{A}\mathbf{X} - \mathbf{A}\mathbb{E}(\mathbf{X}))(\mathbf{B}\mathbf{Y} - \mathbf{B}\mathbb{E}(\mathbf{Y}))^\top\} \\ &= \mathbb{E}\{\mathbf{A}(\mathbf{X} - \mathbb{E}(\mathbf{X}))(\mathbf{Y} - \mathbb{E}(\mathbf{Y}))^\top \mathbf{B}^\top\} \\ &= \mathbf{A} \mathbb{E}\{(\mathbf{X} - \mathbb{E}(\mathbf{X}))(\mathbf{Y} - \mathbb{E}(\mathbf{Y}))^\top\} \mathbf{B}^\top \\ &= \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{Y}) \mathbf{B}^\top.\end{aligned}$$

□

Tenemos el siguiente caso particular,

$$\text{Cov}(\mathbf{A}\mathbf{X}) = \text{Cov}(\mathbf{A}\mathbf{X}, \mathbf{A}\mathbf{X}) = \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{X}) \mathbf{A}^\top = \mathbf{A} \text{Cov}(\mathbf{X}) \mathbf{A}^\top.$$

EJEMPLO 1.7. Considere \mathbf{X} , \mathbf{Y} , \mathbf{U} y \mathbf{V} vectores aleatorios n -dimensionales y \mathbf{A} , \mathbf{B} , \mathbf{C} y \mathbf{D} matrices de órdenes apropiados, entonces

$$\begin{aligned}\text{Cov}(\mathbf{AX} + \mathbf{BY}, \mathbf{CU} + \mathbf{DV}) &= \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{U}) \mathbf{C}^\top + \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{V}) \mathbf{D}^\top \\ &\quad + \mathbf{B} \text{Cov}(\mathbf{Y}, \mathbf{U}) \mathbf{C}^\top + \mathbf{B} \text{Cov}(\mathbf{Y}, \mathbf{V}) \mathbf{D}^\top.\end{aligned}$$

tomando $\mathbf{U} = \mathbf{X}$, $\mathbf{V} = \mathbf{Y}$, $\mathbf{C} = \mathbf{A}$ y $\mathbf{D} = \mathbf{B}$, tenemos

$$\begin{aligned}\text{Cov}(\mathbf{AX} + \mathbf{BY}) &= \text{Cov}(\mathbf{AX} + \mathbf{BY}, \mathbf{AX} + \mathbf{BY}) \\ &= \mathbf{A} \text{Cov}(\mathbf{X}) \mathbf{A}^\top + \mathbf{A} \text{Cov}(\mathbf{X}, \mathbf{Y}) \mathbf{B}^\top \\ &\quad + \mathbf{B} \text{Cov}(\mathbf{Y}, \mathbf{X}) \mathbf{A}^\top + \mathbf{B} \text{Cov}(\mathbf{Y}) \mathbf{B}^\top.\end{aligned}$$

RESULTADO 1.8. *Toda matriz de dispersión es simétrica y semidefinida positiva*

DEMOSTRACIÓN. La simetría de la matriz de dispersión es obvia. Para mostrar que $\text{Cov}(\mathbf{X})$ es semidefinida positiva, sea $\mathbf{Z} = \mathbf{X} - \mathbf{E}(\mathbf{X})$, y considere la variable aleatoria $Y = \mathbf{a}^\top \mathbf{Z}$, para $\mathbf{a} \in \mathbb{R}^n$ un vector arbitrario. Entonces,

$$\begin{aligned}\mathbf{a}^\top \text{Cov}(\mathbf{X}) \mathbf{a} &= \mathbf{a}^\top \mathbf{E}(\mathbf{X} - \mathbf{E}(\mathbf{X}))(\mathbf{X} - \mathbf{E}(\mathbf{X}))^\top \mathbf{a} \\ &= \mathbf{E}(\mathbf{a}^\top (\mathbf{X} - \mathbf{E}(\mathbf{X}))(\mathbf{X} - \mathbf{E}(\mathbf{X}))^\top \mathbf{a}) \\ &= \mathbf{E}(\mathbf{a}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{a}) = \mathbf{E}(Y^2) \geq 0\end{aligned}$$

y por tanto, $\text{Cov}(\mathbf{X})$ es semidefinida positiva.

Ahora, suponga que $\text{Cov}(\mathbf{X})$ es semidefinida positiva de rango r ($r \leq n$). Luego $\text{Cov}(\mathbf{X}) = \mathbf{B} \mathbf{B}^\top$ donde $\mathbf{B} \in \mathbb{R}^{n \times r}$ de rango r . Sea \mathbf{Y} vector aleatorio r -dimensional con $\mathbf{E}(\mathbf{Y}) = \mathbf{0}$ y $\text{Cov}(\mathbf{Y}) = \mathbf{I}$. Haciendo $\mathbf{X} = \mathbf{B} \mathbf{Y}$, sigue que $\mathbf{E}(\mathbf{X}) = \mathbf{0}$ y

$$\text{Cov}(\mathbf{X}) = \text{Cov}(\mathbf{B} \mathbf{Y}) = \mathbf{B} \text{Cov}(\mathbf{Y}) \mathbf{B}^\top = \mathbf{B} \mathbf{B}^\top.$$

Es decir, corresponde a una matriz de covarianza. \square

RESULTADO 1.9. *Sea \mathbf{X} vector aleatorio n -dimensional y considere la transformación lineal $\mathbf{Y} = \mathbf{A} \mathbf{X} + \mathbf{b}$, donde \mathbf{A} es una matriz de constantes $m \times n$ y \mathbf{b} es vector de constantes $m \times 1$. Entonces*

$$\mathbf{E}(\mathbf{Y}) = \mathbf{A} \mathbf{E}(\mathbf{X}) + \mathbf{b}, \quad \text{Cov}(\mathbf{Y}) = \mathbf{A} \text{Cov}(\mathbf{X}) \mathbf{A}^\top.$$

EJEMPLO 1.10. Sea \mathbf{X} vector aleatorio n -dimensional con media $\mathbf{E}(\mathbf{X}) = \boldsymbol{\mu}$ y matriz de dispersión $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$. Sea

$$\boldsymbol{\Sigma} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$$

la descomposición espectral de $\boldsymbol{\Sigma}$, donde \mathbf{U} es matriz ortogonal y $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\lambda})$, y considere la siguiente transformación

$$\mathbf{Z} = \boldsymbol{\Lambda}^{-1/2} \mathbf{U}^\top (\mathbf{X} - \boldsymbol{\mu})$$

de este modo, obtenemos que

$$\mathbf{E}(\mathbf{Z}) = \mathbf{0} \quad \text{y} \quad \text{Cov}(\mathbf{Z}) = \mathbf{I}.$$

En efecto, la transformación $\mathbf{Z} = \boldsymbol{\Sigma}^{-1/2} (\mathbf{X} - \boldsymbol{\mu})$ también satisface que $\mathbf{E}(\mathbf{Z}) = \mathbf{0}$ y $\text{Cov}(\mathbf{Z}) = \mathbf{I}$.

Suponga que \mathbf{Z} es una matriz aleatoria $n \times p$ cuyas filas son vectores aleatorios independientes $p \times 1$, cada uno con la misma matriz de covarianza $\mathbf{\Sigma}$. Considere la partición

$$\mathbf{Z}^\top = (\mathbf{Z}_1, \dots, \mathbf{Z}_n),$$

donde $\text{Cov}(\mathbf{Z}_i) = \mathbf{\Sigma}$, para $i = 1, \dots, n$. Tenemos que

$$\text{vec}(\mathbf{Z}^\top) = \begin{pmatrix} \mathbf{Z}_1 \\ \vdots \\ \mathbf{Z}_n \end{pmatrix},$$

y dado que todos los \mathbf{Z}_i son independientes con la misma matriz de covarianza, podemos escribir

$$\text{Cov}(\text{vec}(\mathbf{Z}^\top)) = \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{\Sigma} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{\Sigma} \end{pmatrix} = \mathbf{I}_n \otimes \mathbf{\Sigma}.$$

Ahora suponga que llevamos a cabo la transformación lineal $\mathbf{Y} = \mathbf{AZB}$, donde $\mathbf{A} \in \mathbb{R}^{r \times n}$, $\mathbf{B} \in \mathbb{R}^{p \times q}$ son matrices de constantes. Entonces $\mathbf{E}(\mathbf{Y}) = \mathbf{A} \mathbf{E}(\mathbf{Z}) \mathbf{B}$, mientras que

$$\text{vec}(\mathbf{Y}^\top) = (\mathbf{A} \otimes \mathbf{B}^\top) \text{vec}(\mathbf{Z}^\top),$$

de modo que

$$\mathbf{E}(\text{vec}(\mathbf{Y}^\top)) = (\mathbf{A} \otimes \mathbf{B}^\top) \mathbf{E}(\text{vec}(\mathbf{Z}^\top)).$$

Lo que lleva a calcular fácilmente la matriz de covarianza

$$\begin{aligned} \text{Cov}(\text{vec}(\mathbf{Y}^\top)) &= (\mathbf{A} \otimes \mathbf{B}^\top) \text{Cov}(\text{vec}(\mathbf{Z}^\top)) (\mathbf{A} \otimes \mathbf{B}^\top)^\top \\ &= (\mathbf{A} \otimes \mathbf{B}^\top) (\mathbf{I}_n \otimes \mathbf{\Sigma}) (\mathbf{A}^\top \otimes \mathbf{B}) \\ &= \mathbf{A} \mathbf{A}^\top \otimes \mathbf{B}^\top \mathbf{\Sigma} \mathbf{B}. \end{aligned}$$

DEFINICIÓN 1.11 (Matriz de correlación). Sea $\mathbf{X} = (X_1, \dots, X_p)^\top$ vector aleatorio con media $\boldsymbol{\mu}$ y matriz de covarianza $\mathbf{\Sigma}$. Se define la matriz de correlaciones como $\mathbf{R} = (\rho_{ij})$, donde

$$\rho_{ij} = \frac{\text{Cov}(X_i, X_j)}{\{\text{var}(X_i) \text{var}(X_j)\}^{1/2}} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii} \sigma_{jj}}}, \quad i, j = 1, \dots, p.$$

Note que, para $\mathbf{\Sigma}$ matriz de covarianza del vector aleatorio \mathbf{X} y con $\mathbf{D} = \text{diag}(\sigma_{11}, \dots, \sigma_{pp})$ ($= \text{diag}(\mathbf{\Sigma})$) podemos escribir

$$\mathbf{R} = \mathbf{D}^{-1/2} \mathbf{\Sigma} \mathbf{D}^{-1/2}.$$

Cada elemento de la diagonal de \mathbf{R} es igual a 1, mientras que sus elementos fuera de la diagonal están entre -1 y 1 . Además se desprende desde la definición que \mathbf{R} es una matriz semidefinida positiva.

RESULTADO 1.12. Sea \mathbf{X} vector aleatorio p -dimensional con $\mathbf{E}(\mathbf{X}) = \boldsymbol{\mu}$ y $\text{Cov}(\mathbf{X}) = \mathbf{\Sigma}$. Sea \mathbf{A} una matriz $p \times p$. Entonces

$$\mathbf{E}(\mathbf{X}^\top \mathbf{A} \mathbf{X}) = \text{tr}(\mathbf{A} \mathbf{\Sigma}) + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}.$$

DEMOSTRACIÓN. Tenemos

$$\begin{aligned}\mathbf{E}(\mathbf{X}^\top \mathbf{A} \mathbf{X}) &= \mathbf{E}(\text{tr } \mathbf{X}^\top \mathbf{A} \mathbf{X}) = \mathbf{E}(\text{tr } \mathbf{A} \mathbf{X} \mathbf{X}^\top) \\ &= \text{tr } \mathbf{E}(\mathbf{A} \mathbf{X} \mathbf{X}^\top) = \text{tr } \mathbf{A} \mathbf{E}(\mathbf{X} \mathbf{X}^\top) \\ &= \text{tr } \mathbf{A}(\boldsymbol{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^\top) = \text{tr}(\mathbf{A} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}.\end{aligned}$$

□

Considere el siguiente caso especial: sea $\mathbf{Y} = \mathbf{X} - \mathbf{a}$, entonces $\text{Cov}(\mathbf{Y}) = \text{Cov}(\mathbf{X})$ y tenemos

$$\mathbf{E}[(\mathbf{X} - \mathbf{a})^\top \mathbf{A}(\mathbf{X} - \mathbf{a})] = \text{tr}(\mathbf{A} \boldsymbol{\Sigma}) + (\boldsymbol{\mu} - \mathbf{a})^\top \mathbf{A}(\boldsymbol{\mu} - \mathbf{a}).$$

EJEMPLO 1.13. Sea $\mathbf{1}_n = (1, \dots, 1)^\top$ vector n -dimensional cuyos componentes son todos 1. Note que, $\mathbf{1}_n^\top \mathbf{1}_n = n$. Considere el vector aleatorio $\mathbf{X} = (X_1, \dots, X_n)^\top$, entonces

$$\mathbf{X}^\top \mathbf{X} = \sum_{i=1}^n X_i^2, \quad \mathbf{1}^\top \mathbf{X} = \sum_{i=1}^n X_i.$$

De este modo, tenemos

$$\begin{aligned}\sum_{i=1}^n (X_i - \bar{X})^2 &= \sum_{i=1}^n X_i^2 - n\bar{X}^2 = \mathbf{X}^\top \mathbf{X} - n\left(\frac{1}{n} \mathbf{1}^\top \mathbf{X}\right)^2 \\ &= \mathbf{X}^\top \mathbf{X} - n\left(\frac{1}{n} \mathbf{1}^\top \mathbf{X}\right)\left(\frac{1}{n} \mathbf{1}^\top \mathbf{X}\right) = \mathbf{X}^\top \mathbf{X} - \frac{1}{n} \mathbf{X}^\top \mathbf{1} \mathbf{1}^\top \mathbf{X} \\ &= \mathbf{X}^\top \left(\mathbf{I} - \frac{1}{n} \mathbf{J}_n\right) \mathbf{X}, \quad \mathbf{J}_n = \mathbf{1}_n \mathbf{1}_n^\top\end{aligned}$$

Llamaremos a $\mathbf{C} = \mathbf{I} - \frac{1}{n} \mathbf{J}_n$ la *matriz de centrado*. Suponga que X_1, \dots, X_n son variables aleatorias independientes e idénticamente distribuidas con media μ y varianza σ^2 . Sigue que,

$$\mathbf{E}(\mathbf{X}) = \mu \mathbf{1}_n, \quad \boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_n,$$

pues $\text{Cov}(X_i, X_j) = 0$ ($i \neq j$). Por tanto, podemos usar el Resultado (1.12) para calcular la esperanza de la variable aleatoria,

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

obteniendo

$$\mathbf{E}(Q) = \sigma^2 \text{tr}(\mathbf{C}) + \mu^2 \mathbf{1}^\top \mathbf{C} \mathbf{1}.$$

Es fácil verificar que

$$\begin{aligned}\text{tr}(\mathbf{C}) &= \text{tr}\left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top\right) = \text{tr}(\mathbf{I}) - \frac{1}{n} \text{tr}(\mathbf{1} \mathbf{1}^\top) = n - \frac{1}{n} \mathbf{1}^\top \mathbf{1} = n - 1, \\ \mathbf{C} \mathbf{1} &= \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top\right) \mathbf{1} = \mathbf{1} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \mathbf{1} = \mathbf{1} - \mathbf{1} = \mathbf{0},\end{aligned}$$

de donde sigue que $\mathbf{E}(Q) = \sigma^2(n - 1)$.

RESULTADO 1.14. Si \mathbf{X} es vector aleatorio $n \times 1$. Entonces su distribución está determinada por las distribuciones de las funciones lineales $\mathbf{a}^\top \mathbf{X}$, para todo $\mathbf{a} \in \mathbb{R}^n$.

DEMOSTRACIÓN. La función característica de $\mathbf{a}^\top \mathbf{X}$ es

$$\varphi_{\mathbf{a}^\top \mathbf{X}}(t) = \mathbf{E}\{\exp(it\mathbf{a}^\top \mathbf{X})\},$$

de modo que

$$\varphi_{\mathbf{a}^\top \mathbf{X}}(1) = \mathbf{E}\{\exp(i\mathbf{a}^\top \mathbf{X})\} = \varphi_X(\mathbf{a}).$$

Es considerada como una función de \mathbf{a} , esto es, la función característica (conjunta) de \mathbf{X} . El resultado sigue notando que una distribución en \mathbb{R}^n está completamente determinada por su función característica. \square

La función característica permite un método bastante operativo para el cálculo del k -ésimo momento de un vector aleatorio \mathbf{X} . En efecto,

$$\begin{aligned} \mu_k(\mathbf{X}) &= \begin{cases} \mathbf{E}(\mathbf{X} \otimes \mathbf{X}^\top \otimes \cdots \otimes \mathbf{X}^\top), & k \text{ par,} \\ \mathbf{E}(\mathbf{X} \otimes \mathbf{X}^\top \otimes \cdots \otimes \mathbf{X}^\top \otimes \mathbf{X}), & k \text{ impar,} \end{cases} \\ &= \begin{cases} i^{-k} \frac{\partial^k \varphi(\mathbf{t})}{\partial \mathbf{t} \partial \mathbf{t}^\top \cdots \partial \mathbf{t}^\top} \Big|_{\mathbf{t}=\mathbf{0}}, & k \text{ par,} \\ i^{-k} \frac{\partial^k \varphi(\mathbf{t})}{\partial \mathbf{t} \partial \mathbf{t}^\top \cdots \partial \mathbf{t}^\top \partial \mathbf{t}} \Big|_{\mathbf{t}=\mathbf{0}}, & k \text{ impar.} \end{cases} \end{aligned}$$

1.3. Independencia de vectores aleatorios

Sea $\mathbf{Z} = (\mathbf{X}^\top, \mathbf{Y}^\top)^\top$ con \mathbf{X} , \mathbf{Y} vectores aleatorios n y q -dimensionales, respectivamente. Se dicen independientes si y sólo si

$$F(\mathbf{x}, \mathbf{y}) = G(\mathbf{x})H(\mathbf{y}),$$

donde $F(\mathbf{z})$, $G(\mathbf{x})$ y $H(\mathbf{y})$ son las funciones de distribución de \mathbf{Z} , \mathbf{X} e \mathbf{Y} , respectivamente.

Si \mathbf{Z} , \mathbf{X} e \mathbf{Y} tienen densidades $f(\mathbf{z})$, $g(\mathbf{x})$ y $h(\mathbf{y})$, respectivamente. Entonces \mathbf{X} e \mathbf{Y} son independientes si

$$f(\mathbf{z}) = g(\mathbf{x})h(\mathbf{y}).$$

En cuyo caso, obtenemos como resultado

$$f(\mathbf{x}|\mathbf{y}) = g(\mathbf{x}).$$

RESULTADO 1.15. Sean \mathbf{X} e \mathbf{Y} dos vectores aleatorios independientes. Entonces para funciones cualquiera κ y τ , tenemos

$$\mathbf{E}\{\kappa(\mathbf{X})\tau(\mathbf{Y})\} = \mathbf{E}\{\kappa(\mathbf{X})\}\mathbf{E}\{\tau(\mathbf{Y})\},$$

si las esperanzas existen.

DEMOSTRACIÓN. En efecto, es fácil notar que

$$\begin{aligned} \mathbf{E}\{\kappa(\mathbf{X})\tau(\mathbf{Y})\} &= \int \int \kappa(\mathbf{x})\tau(\mathbf{y})g(\mathbf{x})h(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} \\ &= \left(\int \kappa(\mathbf{x})g(\mathbf{x}) \, d\mathbf{x} \right) \left(\int \tau(\mathbf{y})h(\mathbf{y}) \, d\mathbf{y} \right) \\ &= \mathbf{E}\{\kappa(\mathbf{X})\}\mathbf{E}\{\tau(\mathbf{Y})\}. \end{aligned}$$

\square

1.4. Cambios de variable

Considere la función $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, el *Jacobiano* se define como el valor absoluto del determinante de $D\mathbf{f}(\mathbf{x})$ y es denotado por

$$J(\mathbf{y} \rightarrow \mathbf{x}) = |D\mathbf{f}(\mathbf{x})|_+ = \text{abs}(\det(D\mathbf{f}(\mathbf{x}))),$$

donde $\mathbf{y} = \mathbf{f}(\mathbf{x})$. Note que si $\mathbf{z} = \mathbf{f}(\mathbf{y})$ y $\mathbf{y} = \mathbf{g}(\mathbf{x})$, entonces tenemos

$$J(\mathbf{z} \rightarrow \mathbf{x}) = J(\mathbf{z} \rightarrow \mathbf{y}) \cdot J(\mathbf{y} \rightarrow \mathbf{x})$$

$$J(\mathbf{y} \rightarrow \mathbf{x}) = \{J(\mathbf{x} \rightarrow \mathbf{y})\}^{-1}$$

El siguiente resultado presenta una de aplicación del Jacobiano de una transformación para obtener la función de densidad de una transformación de un vector aleatorio.

PROPOSICIÓN 1.16 (Transformación de vectores aleatorios continuos). *Sea \mathbf{X} vector aleatorio n -dimensional con densidad $f_X(\mathbf{x})$ y soporte $S = \{\mathbf{x} : f_X(\mathbf{x}) > 0\}$. Para $\mathbf{g} : S \rightarrow \mathbb{R}^n$ diferenciable e invertible, sea $\mathbf{y} = \mathbf{g}(\mathbf{x})$. Entonces la densidad de \mathbf{Y} está dada por*

$$\begin{aligned} f_Y(\mathbf{y}) &= |D\mathbf{g}^{-1}(\mathbf{y})|_+ f_X(\mathbf{g}^{-1}(\mathbf{y})) \\ &= \{J(\mathbf{y} \rightarrow \mathbf{x})\}^{-1} f_X(\mathbf{g}^{-1}(\mathbf{y})). \end{aligned}$$

EJEMPLO 1.17. Sea $\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{B}$, $\mathbf{Y} \in \mathbb{R}^{n \times q}$, $\mathbf{X} \in \mathbb{R}^{n \times q}$, $\mathbf{A} \in \mathbb{R}^{n \times n}$ y $\mathbf{B} \in \mathbb{R}^{q \times q}$. Entonces

$$d\mathbf{Y} = \mathbf{A}(d\mathbf{X})\mathbf{B},$$

vectorizando obtenemos

$$\text{vec } d\mathbf{Y} = (\mathbf{B}^\top \otimes \mathbf{A}) \text{vec } d\mathbf{X},$$

esto es, $D\mathbf{F}(\mathbf{X}) = \mathbf{B}^\top \otimes \mathbf{A}$, por tanto

$$J(\mathbf{Y} \rightarrow \mathbf{X}) = |\mathbf{B}^\top \otimes \mathbf{A}|_+ = |\mathbf{A}|_+^q |\mathbf{B}^\top|_+^n = |\mathbf{A}|_+^q |\mathbf{B}|_+^n$$

1.5. Modelo estadístico

El punto de partida es considerar que los datos observados \mathbf{y} son una realización de una variable aleatoria Y cuya distribución es parcialmente conocida. Asumiremos que la distribución de la variable de interés es un miembro de una clase de medidas de probabilidad definida como

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

que es indexada por $\theta \in \Theta$. El conjunto Θ es denominado *espacio paramétrico*.

IDEA. Se desea conocer algunas propiedades de \mathcal{P} basado en una (única) muestra aleatoria

DEFINICIÓN 1.18. Un *modelo estadístico* es un par $(\mathcal{Y}, \mathcal{P})$, donde \mathcal{Y} es el conjunto de todos los resultados posibles, conocido como *espacio muestral*.

DEFINICIÓN 1.19. Una *muestra aleatoria*¹ $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ es una colección de variables aleatorias independientes, donde cada Y_i es distribuída como P_θ . En este contexto, n corresponde al número de variables aleatorias y es llamado *tamaño muestral*.

¹Usualmente utilizaremos la notación $\text{m.a.}(n)$.

En la práctica, se dispondrá de observaciones que son realizaciones de variables aleatorias independientes Y_1, \dots, Y_n . En este caso \mathcal{P} es una familia de medidas producto

$$\mathbf{P}_\theta = \mathbf{P}_{1,\theta} \otimes \cdots \otimes \mathbf{P}_{n,\theta}.$$

Una muestra aleatoria IID, $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$, es una colección de variables aleatorias independientes, donde cada Y_i tiene la misma distribución. Es decir, $\mathbf{P}_{i,\theta} = \mathbf{P}_{1,\theta}$ para todo i .

Evidentemente, cada Y_i tendrá función de distribución acumulada *común*, digamos F . Si F es *conocido*, podemos usar cálculo de probabilidades para deducir y estudiar sus propiedades. Sin embargo, en la práctica, F es *desconocido*, y el objetivo es tratar de inferir sus propiedades *desde los datos*. Frecuentemente, el interés es una *función no aleatoria* de F , tal como la media o su q -ésimo cuantil,

$$\mathbf{E}(Y) = \int y \, dF(y), \quad y_q = F^{-1}(q) = \inf\{y : F(y) \geq q\}.$$

Las cantidades $\mathbf{E}(Y)$, $\text{var}(Y)$ y $F^{-1}(q)$ son llamadas *parámetros*, dependen de F y son generalmente desconocidos.

OBSERVACIÓN. La familia $\mathcal{P} = \{\mathbf{P}_\theta : \theta \in \Theta \subset \mathbb{R}^k\}$ es parametrizada y decimos que el modelo es *paramétrico*.²

OBSERVACIÓN. Recuerde que, para $A \subset \mathcal{Y}$, la función de probabilidad para el caso continuo es dada por

$$\mathbf{P}_\theta(A) = \int_A f(y; \theta) \, dy,$$

donde $f(y; \theta)$ es la función de densidad de Y . Mientras que para el caso discreto, tenemos

$$\mathbf{P}_\theta(A) = \sum_{y \in A} \mathbf{P}_\theta(\{y\}).$$

En efecto, usaremos la notación $\mathbf{P}_\theta(\mathbf{y})$ o $\mathbf{P}_\theta(\mathbf{Y} = \mathbf{y})$ en lugar de $\mathbf{P}_\theta(\{\mathbf{y}\})$.

EJEMPLO 1.20. Considere $X_i \sim \text{Bin}(n, \theta)$ es decir

$$p(x_i; \theta) = \binom{n}{x_i} \theta^{x_i} (1 - \theta)^{n-x_i},$$

con $n \in \mathbb{N}$, $\theta \in (0, 1)$ y $x_i \in \{0, 1, \dots, n\}$. Es decir $\mathcal{X} = \{0, 1, \dots, n\}$ y $\Theta = (0, 1)$ (estamos considerando n fijo). De este modo,

$$\mathcal{P} = \{\text{Bin}(n, \theta) : \theta \in (0, 1)\}.$$

EJEMPLO 1.21. Suponga que $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ son IID tal que $Y_i \sim \text{Poi}(\theta)$, para $i = 1, \dots, n$. Aquí

$$p(y_i; \theta) = \frac{e^{-\theta} \theta^{y_i}}{y_i!}, \quad \theta \in (0, \infty),$$

y $y_i \in \{0, 1, \dots\}$. Esto lleva al modelo estadístico

$$\mathcal{P} = \{\text{Poi}(\theta)^{\otimes n} : \theta \in (0, \infty)\}.$$

²Si k no es entero, decimos que el modelo es *no paramétrico*.

Note que la *densidad conjunta* (asociada a $\text{Poi}(\theta)^{\otimes n}$) es dada por

$$p(\mathbf{y}; \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{y_i}}{y_i!} = \frac{e^{-n\theta} \theta^{\sum_{i=1}^n y_i}}{\prod_{i=1}^n y_i!}.$$

EJEMPLO 1.22. Suponga $\mathbf{Y} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ con

$$f(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \right\}.$$

En este caso tenemos, $\mathcal{Y} = \mathbb{R}^p$, mientras que el modelo estadístico es definido como

$$\mathcal{P} = \{\mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) : \boldsymbol{\mu} \in \mathbb{R}^p, \boldsymbol{\Sigma} > 0\}.$$

EJEMPLO 1.23. Se ha sugerido que el crecimiento (durante siete semanas) de plantas de soya puede ser descrito por el modelo

$$Y_t = a + bt + \epsilon_t, \quad t = 1, \dots, 7.$$

Además, asumiremos que los errores satisfacen

$$\mathbf{E}(\epsilon_t) = 0, \quad \text{var}(\epsilon_t) = \sigma^2,$$

y $\mathbf{E}(\epsilon_t \epsilon_r) = 0$ para $t \neq r$. De esta manera el modelo para $\mathbf{Y} = (Y_1, \dots, Y_7)^\top$ toma valores en $\mathcal{Y} = \mathbb{R}^7$ y puede ser escrito como:

$$\begin{aligned} \mathcal{P} &= \{\mathbf{P}_\theta : \mathbf{E}(Y_t) = a + bt, \text{var}(Y_t) = \sigma^2, \quad 1 \leq t \leq 7; \\ &\quad \boldsymbol{\theta} = (a, b, \sigma^2) \in \Theta \subset \mathbb{R}^2 \times \mathbb{R}_+\}, \end{aligned}$$

Usualmente en este tipo de problemas se asume que la distribución que caracteriza los primeros momentos de $\{\epsilon_t\}$ es conocida.

OBSERVACIÓN (Parametrización). La parametrización usada no es única. En lugar de $\mathcal{P} = \{\mathbf{P}_\theta : \theta \in \Theta\}$, podemos escoger

$$\mathcal{P} = \{\mathbf{P}_\phi : \phi \in \Phi\},$$

donde $\theta = h^{-1}(\phi)$ para h función 1 a 1.

EJEMPLO 1.24. Considere $X \sim \text{Exp}(\lambda)$ con densidad

$$f(x; \lambda) = \lambda \exp(-\lambda x), \quad x \geq 0, \lambda > 0.$$

Tenemos la parametrización alternativa, considerando $\mathbf{E}(X) = 1/\lambda = \mu$. De este modo,

$$f(x; \mu) = \frac{1}{\mu} \exp(-x/\mu), \quad x \geq 0, \mu > 0.$$

SUPUESTO A0 (Identificabilidad). Para $\theta_1, \theta_2 \in \Theta$ con $\theta_1 \neq \theta_2$. Si las distribuciones \mathbf{P}_{θ_1} y \mathbf{P}_{θ_2} son diferentes, entonces se dice que el modelo es **identificable**.

La elección del modelo estadístico es esencial para realizar inferencia. Aunque todos los modelos están errados, algunos resultarán útiles y se deberá mantener el modelo tan sencillo como sea posible. Por tanto, una tarea crucial será evaluar la validez del modelo propuesto.

DEFINICIÓN 1.25. Una *estadística* T es una función de la muestra. Es decir,

$$T : \mathcal{Y} \rightarrow \mathcal{T},$$

tal que $T(\mathbf{y}) = t$. Note además que $T(\mathbf{Y})$ es una variable aleatoria con función de densidad $f(t; \theta)$.

EJEMPLO 1.26. Sean Y_1, \dots, Y_n variables aleatorias IID desde $N(\mu, \sigma^2)$. En este caso,

$$\mathcal{P} = \{P_\theta : \theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+\}.$$

Ahora, considere

$$T_1(\mathbf{Y}) = \frac{1}{n} \sum_{i=1}^n Y_i, \quad T_2(\mathbf{Y}) = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

Tenemos $\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^2$, aquí $\mathbf{T}(Y_1, \dots, Y_n) \rightarrow (T_1, T_2)$.

EJEMPLO 1.27. Considere la función de distribución empírica

$$\hat{F}(X_1, \dots, X_n)(x) = \hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x),$$

donde $(X_1, \dots, X_n)^\top$ es una muestra aleatoria desde la distribución F (asociado al modelo \mathcal{P}) y $I(A)$ representa la función indicadora del evento A . Es decir, tenemos que $\hat{F} : \mathbb{R}^n \rightarrow \mathbb{R}$ es una estadística.

OBSERVACIÓN. Un estadístico depende *sólo* de la muestra. **No** puede depender de cantidades desconocidas. Por ejemplo, para poblaciones normales

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i (= T_1), \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 (= T_2),$$

son estadísticos para μ y σ^2 , respectivamente. Mientras que

$$T_3 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2,$$

no es un estadístico para σ^2 pues depende de μ .

1.6. Familia exponencial

DEFINICIÓN 1.28. La clase de modelos $\{P_\theta : \theta \in \Theta\}$ se dice una *familia exponencial* 1-paramétrica si existen funciones $\eta(\theta)$, $b(\theta)$ y funciones real valuadas T y h tal que su densidad adopta la forma

$$f(x) = \exp[\eta(\theta)T(x) - b(\theta)]h(x), \quad (1.3)$$

donde $x \in \mathcal{X} \subset \mathbb{R}^q$. Cuando una variable aleatoria tiene la densidad en Ecuación (1.3), anotamos $X \sim \text{FE}(\theta)$.

OBSERVACIÓN. El parámetro $\eta := \eta(\theta)$ es llamado *parámetro natural*, y

$$b(\eta) = \log \int \exp(\eta T(x)) h(x) \, dx,$$

corresponde a un factor de normalización. Además, el modelo estadístico puede ser escrito como:

$$\mathcal{P} = \{P_\eta : \eta \in \Gamma\}, \quad \Gamma := \eta(\Theta), \quad (1.4)$$

donde el conjunto

$$\Gamma = \{\eta : b(\eta) < \infty\},$$

se denomina *espacio paramétrico natural*. Cuando escribimos la familia de densidades $\{P_\eta : \eta \in \Gamma\}$ como en (1.4), decimos que la familia está escrita en *forma canónica*.

EJEMPLO 1.29. Sea $P_\theta \stackrel{d}{=} \text{Poi}(\theta)$ con parámetro de media desconocido. Entonces $x \in \{0, 1, 2, \dots\}$, y

$$p(x; \theta) = \frac{\theta^x e^{-\theta}}{x!} = \frac{1}{x!} \exp(x \log \theta - \theta), \quad \theta > 0.$$

Es decir, P_θ corresponde a una FE (1-paramétrica) con $q = 1$, $\eta(\theta) = \log \theta$, $T(x) = x$, $b(\theta) = \theta$ y $h(x) = 1/x!$.

EJEMPLO 1.30. Sea $P_\theta \stackrel{d}{=} \text{Bin}(n, \theta)$, $\theta \in (0, 1)$ y $x \in \{0, 1, \dots, n\}$. De este modo,

$$\begin{aligned} p(x; \theta) &= \binom{n}{x} \theta^x (1 - \theta)^{n-x} \\ &= \binom{n}{x} \exp[x \log \theta + (n - x) \log(1 - \theta)] \\ &= \binom{n}{x} \exp \left[x \log \left(\frac{\theta}{1 - \theta} \right) + n \log(1 - \theta) \right], \end{aligned}$$

y $\text{Bin}(n, \theta)$ está en FE (1-paramétrica) con $q = 1$, $\eta(\theta) = \log(\theta/(1 - \theta))$, $T(x) = x$, $b(\theta) = -n \log(1 - \theta)$ y $h(x) = \binom{n}{x}$.

OBSERVACIÓN. Note que las funciones η , b y T *no son únicas*. En efecto, una manera alternativa de escribir la densidad de la familia exponencial es, por ejemplo,

$$f(x) = a(\theta) \exp[\eta(\theta)T(x)]h(x),$$

donde $h(x) \geq 0$ y $a(\theta) > 0$, representa un factor de normalización.

EJEMPLO 1.31 (contraejemplo). Sea \mathcal{P} la clase de distribuciones exponencial de dos parámetros, digamos $\text{Exp}(1, \theta)$, esto es con densidad

$$f(x) = \exp[-(x - \theta)]I_{[\theta, \infty)}(x), \quad \theta \in \mathbb{R}.$$

Esta familia no pertenece a la clase FE.³

Suponga ahora que X_1, \dots, X_n son variables aleatorias IID con distribución común P_θ en la FE 1-paramétrica y $\theta \in \Theta$. De este modo,

$$\begin{aligned} f(\mathbf{x}) &= \prod_{i=1}^n \exp[\eta(\theta)T(x_i) - b(\theta)]h(x_i) \\ &= \exp \left[\eta(\theta) \sum_{i=1}^n T(x_i) - nb(\theta) \right] \prod_{i=1}^n h(x_i) \\ &= \exp[\eta(\theta)T^{(n)}(\mathbf{x}) - b^{(n)}(\theta)]\tilde{h}(\mathbf{x}) \end{aligned}$$

con

$$T^{(n)}(\mathbf{x}) = \sum_{i=1}^n T(x_i), \quad b^{(n)}(\theta) = nb(\theta), \quad \tilde{h}(\mathbf{x}) = \prod_{i=1}^n h(x_i). \quad (1.5)$$

Es decir, la distribución conjunta de $\mathbf{X} = (X_1, \dots, X_n)^\top$ también pertenece a la FE 1-paramétrica.

³El factor $I_{[\theta, \infty)}(x)$ no puede ser escrito en la forma exponencial.

RESULTADO 1.32. Si X está en la FE 1-paramétrica. La función generadora de momentos de $T(X)$ es dada por

$$M_T(s) = \exp[b(s + \eta) - b(\eta)],$$

para s en una vecindad de cero.

DEMOSTRACIÓN. Para obtener la MGF de $T(X)$, debemos calcular

$$\begin{aligned} M_T(s) &= \mathbf{E}\{\exp(sT(x))\} = \int_{\mathcal{X}} \exp(sT(x)) \exp(\eta T(x) - b(\eta)) h(x) \, dx \\ &= \int_{\mathcal{X}} \exp[(s + \eta)T(x) - b(\eta)] h(x) \, dx \\ &= \exp[b(s + \eta) - b(\eta)] \int_{\mathcal{X}} \exp[(s + \eta)T(x) - b(s + \eta)] h(x) \, dx \\ &= \exp[b(s + \eta) - b(\eta)]. \end{aligned}$$

□

La función generadora de cumulantes está definida como

$$K_T(s) = \log M_T(s),$$

y permite obtener los cumulantes de T ,

$$\kappa_r = \left. \frac{d^r}{ds^r} K_T(s) \right|_{s=0}.$$

En particular las primeras dos derivadas de $K_T(s)$ asumen la forma,

$$K'_T = M'_T/M_T, \quad K''_T = [M_T M''_T - (M'_T)^2]/M_T^2,$$

y evaluando en $s = 0$, lleva a los primeros cumulantes:

$$\kappa_1 = \mathbf{E}(T), \quad \kappa_2 = \mathbf{E}(T^2) - (\mathbf{E}(T))^2 = \text{var}(T).$$

Para la familia exponencial, sigue que la función generadora de cumulantes para T es dada por:

$$K_T(s) = b(s + \eta) - b(\eta).$$

RESULTADO 1.33. Sea X una variable aleatoria perteneciente a la FE 1-paramétrica. De este modo,

$$\mathbf{E}(T) = b'(\eta), \quad \text{var}(T) = b''(\eta).$$

DEMOSTRACIÓN. Basta notar que

$$\begin{aligned} \mathbf{E}(T) &= \left. \frac{d}{ds} K_T(s) \right|_{s=0} = b'(s + \eta) \Big|_{s=0} = b'(\eta), \\ \text{var}(T) &= \left. \frac{d^2}{ds^2} K_T(s) \right|_{s=0} = b''(s + \eta) \Big|_{s=0} = b''(\eta). \end{aligned}$$

□

Para una muestra IID, X_1, \dots, X_n , y usando el Resultado 1.33 sigue inmediatamente que (ver Ecuación (1.5))

$$\begin{aligned} \mathbb{E}(T^{(n)}(\mathbf{X})) &= \mathbb{E}\left(\sum_{i=1}^n T(X_i)\right) = \sum_{i=1}^n \mathbb{E}(T(X_i)) = nb'(\theta) \\ \text{var}(T^{(n)}(\mathbf{X})) &= \text{var}\left(\sum_{i=1}^n T(X_i)\right) = \sum_{i=1}^n \text{var}(T(X_i)) = nb''(\theta). \end{aligned}$$

Sea $\mu = \mathbb{E}(Y)$, es decir, tenemos una transformación⁴ 1-1 entre θ y el parámetro de media

$$\mu = \mu(\theta) = b'(\theta),$$

pues $\eta = \eta(\theta)$. Además, conforme θ varia en Θ , tenemos que $\mu(\theta)$ es conocida como *función (o superficie) de esperanza*.

Cuando $Y \sim \text{FE}(\theta)$ con función generadora de cumulantes $b(\cdot)$, tenemos $\mu(\theta) = b'(\theta)$. Podemos parametrizar la densidad $f(y)$ en términos de μ . De este modo, haciendo $\theta = \theta(\mu)$, sigue que

$$\text{var}(Y) = b''(\theta) = \left. \frac{d\mu}{d\theta} \right|_{\theta=\theta(\mu)} = V(\mu),$$

donde $V(\mu)$ es llamada *función de varianza* de Y .

DEFINICIÓN 1.34. La familia de distribuciones $\{P_\theta : \theta \in \Theta\}$ con $\Theta \subset \mathbb{R}^k$, se dice una *familia exponencial k-paramétrica* si existen funciones $\eta_1(\theta), \dots, \eta_k(\theta)$ y $b(\theta)$ y funciones real-valuadas T_1, \dots, T_k , h tal que la densidad puede ser escrita como

$$f(\mathbf{x}; \theta) = \exp[\boldsymbol{\eta}^\top(\theta) \mathbf{T}(\mathbf{x}) - b(\theta)] h(\mathbf{x}), \quad \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^p,$$

donde $\boldsymbol{\eta}(\theta) = (\eta_1(\theta), \dots, \eta_k(\theta))^\top$ y $\mathbf{T}(\mathbf{x}) = (T_1(\mathbf{x}), \dots, T_k(\mathbf{x}))^\top$.

EJEMPLO 1.35. Suponga $P_\theta \stackrel{d}{=} N_1(\mu, \sigma^2)$,

$$\Theta = \{(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}.$$

La densidad de P_θ puede ser escrita como

$$\begin{aligned} f(x; \theta) &= (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\} \\ &= \exp\left\{\frac{\mu x}{\sigma^2} - \frac{x^2}{2\sigma^2} - \frac{1}{2}\left(\frac{\mu^2}{\sigma^2} + \log 2\pi\sigma^2\right)\right\}. \end{aligned}$$

De este modo,

$$\eta_1(\theta) = \frac{\mu}{\sigma^2}, \quad \eta_2(\theta) = -\frac{1}{2\sigma^2}, \quad T_1(x) = x, \quad T_2(x) = x^2,$$

y $b(\theta) = \frac{1}{2}(\mu^2/\sigma^2 + \log 2\pi\sigma^2)$, $h(x) = 1$. De este modo, $N_1(\mu, \sigma)$ corresponde a una FE 2-paramétrica.

EJEMPLO 1.36. Considere la función de probabilidad

$$p(x; \theta) = \theta_1^{I_1(x)} \theta_2^{I_2(x)} \theta_3^{I_3(x)},$$

con

$$I_j(x) = \begin{cases} 1, & x = j, \\ 0, & \text{en otro caso,} \end{cases}$$

⁴suave y estrictamente convexa

para $j = 1, 2, 3$. Así podemos escribir

$$p(x; \boldsymbol{\theta}) = \exp[I_1(x) \log \theta_1 + I_2(x) \log \theta_2 + I_3(x) \log \theta_3],$$

es decir X pertenece a la FE. Sin embargo,

$$I_1(x) + I_2(x) + I_3(x) = 1.$$

Por tanto, la familia **no** es estrictamente 3-dimensional. Considere

$$\begin{aligned} p(x; \boldsymbol{\theta}) &= \exp[I_1(x) \log \theta_1 + I_2(x) \log \theta_2 + (1 - I_1(x) - I_2(x)) \log \theta_3] \\ &= \theta_3 \exp[I_1(x)(\log \theta_1 - \log \theta_3) + I_2(x)(\log \theta_2 - \log \theta_3)]. \end{aligned}$$

Sea $\theta_3 = 1 - \theta_1 - \theta_2$, sigue que

$$p(x; \boldsymbol{\theta}) = \exp \left[I_1(x) \log \left(\frac{\theta_1}{1 - \theta_1 - \theta_2} \right) + I_2(x) \log \left(\frac{\theta_2}{1 - \theta_1 - \theta_2} \right) + \log(1 - \theta_1 - \theta_2) \right].$$

Finalmente, tenemos que X pertenece a la FE (2-paramétrica) con

$$T_1(x) = I_1(x), \quad T_2(x) = I_2(x), \quad \eta_r = \log \left(\frac{\theta_r}{1 - \theta_1 - \theta_2} \right), \quad r = 1, 2,$$

$$\text{y } b(\boldsymbol{\theta}) = -\log(1 - \theta_1 - \theta_2), \quad h(x) = 1.$$

OBSERVACIÓN. Evidentemente, resulta fácil extender los Resultados 1.32 y 1.33. En efecto,

$$M_T(\mathbf{s}) = \mathbb{E}[\exp(\mathbf{s}^\top \mathbf{T})] = \exp[b(\boldsymbol{\eta} + \mathbf{s}) - b(\boldsymbol{\eta})], \quad K_T(\mathbf{s}) = b(\boldsymbol{\eta} + \mathbf{s}) - b(\boldsymbol{\eta}),$$

y

$$\mathbb{E}(\mathbf{T}(\mathbf{X})) = \dot{b}(\boldsymbol{\eta}) = \frac{\partial b(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}}, \quad \text{Cov}(\mathbf{T}(\mathbf{X})) = \ddot{b}(\boldsymbol{\eta}) = \frac{\partial^2 b(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^\top}.$$

En aplicación de regresión,⁵ es útil considerar una definición alternativa de la familia exponencial, incluyendo un parámetro de dispersión. Considere la siguiente definición

DEFINICIÓN 1.37. Se dice que $\mathbf{P}_{\theta, \phi}$ sigue una familia exponencial con parámetro de escala $\phi > 0$, si la función de densidad de un vector aleatorio n -dimensional $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ es dada por:

$$f(\mathbf{y}; \boldsymbol{\theta}, \phi) = \exp[\phi \{\mathbf{y}^\top \boldsymbol{\theta} - b(\boldsymbol{\theta})\} + c(\mathbf{y}, \phi)], \quad (1.6)$$

donde $b(\cdot)$ y $c(\cdot, \cdot)$ son funciones apropiadas, $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)^\top$ denota el parámetro natural ($\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^n$).

El parámetro de dispersión $\phi > 0$ usualmente es considerado un parámetro molesto y cuando un vector aleatorio tiene la densidad dada en (1.6), anotamos $\mathbf{Y} \sim \text{FE}(\boldsymbol{\theta}, \phi)$.

OBSERVACIÓN. Cuando ϕ es desconocido, se suele indicar que el modelo está en la familia de dispersión exponencial (Jørgensen, 1997).

EJEMPLO 1.38. Suponga $\mathbf{Y} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Tenemos que su función de densidad puede ser escrita como

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

⁵En particular en modelos lineales generalizados (GLM)

de este modo, el parámetro natural es $\boldsymbol{\theta} = \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}$ y $b(\boldsymbol{\theta}) = \frac{1}{2}\boldsymbol{\theta}^\top \boldsymbol{\Sigma}\boldsymbol{\theta}$, mientras que $c(\mathbf{y}, \boldsymbol{\Phi}) = -\frac{1}{2}(\mathbf{y}^\top \boldsymbol{\Phi} \mathbf{y} - \log |2\pi \boldsymbol{\Phi}^{-1}|)$.

EJEMPLO 1.39. Usando que

$$\mathbf{y}^\top \boldsymbol{\Sigma}^{-1} \mathbf{y} = \text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{y} \mathbf{y}^\top = (\text{vec } \mathbf{y} \mathbf{y}^\top)^\top \text{vec } \boldsymbol{\Sigma}^{-1},$$

podemos escribir \mathbf{Y} como un miembro de la familia exponencial, considerando

$$\mathbf{T}_1(\mathbf{Y}) = \mathbf{Y}, \quad \mathbf{T}_2(\mathbf{Y}) = \mathbf{Y} \mathbf{Y}^\top, \quad \boldsymbol{\eta}_1(\boldsymbol{\theta}) = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}, \quad \boldsymbol{\eta}_2(\boldsymbol{\theta}) = \text{vec } \boldsymbol{\Sigma}^{-1}.$$

Las siguientes clases de distribuciones pueden ser caracterizadas mediante transformaciones de vectores aleatorios.

1.7. Familia de posición-escala

Considere U una variable aleatoria con función de distribución F . De este modo, la variable

$$X = U + a,$$

tiene función de distribución

$$\mathbf{P}(X \leq x) = F(x - a),$$

para F fijo y $a \in \mathbb{R}$ tenemos que X corresponde a la *familia de posición*. Análogamente la *familia de escala* es generada por la transformación

$$X = bU, \quad b > 0,$$

en cuyo caso,

$$\mathbf{P}(X \leq x) = F(x/b), \quad b > 0.$$

Esto lleva a la siguiente definición

DEFINICIÓN 1.40. Sea U una variable aleatoria con función de distribución acumulada fija F y considere la transformación

$$X = a + bU, \quad a \in \mathbb{R}, b > 0.$$

Tenemos

$$\mathbf{P}(X \leq x) = F\left(\frac{x - a}{b}\right).$$

De este modo, X es conocida como *familia de posición-escala*.

Usualmente asociado a F tenemos una función de densidad f , dada por

$$f(x; a, b) = \frac{d}{dx} F\left(\frac{x - a}{b}\right) = \frac{1}{b} F'\left(\frac{x - a}{b}\right) = \frac{1}{b} f\left(\frac{x - a}{b}\right).$$

EJEMPLO 1.41. Algunas familias de posición-escala corresponden a:

- Normal, $\mathbf{N}(a, b^2)$:

$$f(y; a, b) = \frac{1}{b} (2\pi)^{-1/2} \exp \left\{ -\frac{1}{2} \left(\frac{y - a}{b} \right)^2 \right\}.$$

- Laplace (doble exponencial), $\text{Laplace}(a, b)$:

$$f(y; a, b) = \frac{1}{2b} \exp \left\{ -\frac{|y - a|}{b} \right\}.$$

- Cauchy, $\text{Cauchy}(a, b^2)$:

$$f(y; a, b) = \frac{b}{\pi} \frac{1}{b^2 + (y - a)^2}.$$

- Logística, $\text{Logística}(a, b)$:

$$f(y; a, b) = \frac{1}{b} \frac{e^{-(y-a)/b}}{(1 + e^{-(y-a)/b})^2}$$

1.8. Familia de distribuciones de contornos elípticos

Esta clase de distribuciones es ampliamente utilizada en modelamiento estadístico desde la perspectiva de robustez, y corresponde a una extensión natural de la distribución normal multivariada así como de la familia de posición-escala. Para introducir ideas, primeramente revisaremos algunas definiciones de la normal.

DEFINICIÓN 1.42. Un vector aleatorio p -dimensional, \mathbf{X} tiene distribución normal con vector de medias $\boldsymbol{\mu} \in \mathbb{R}^p$ y matriz de covarianza $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma} > 0$ sólo si, para todo vector \mathbf{t} la variable aleatoria (uni-dimensional) $\mathbf{t}^\top \mathbf{X}$ es normal y escribimos $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Note que en la definición anterior **no** se ha hecho supuestos respecto de la independencia de los componentes de \mathbf{X} .

DEFINICIÓN 1.43. La función característica de $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ está dada por

$$\varphi_{\mathbf{X}}(\mathbf{t}) = \exp(i\mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}).$$

Podemos derivar la función característica de $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ desde la normal univariada notando que $Y = \mathbf{t}^\top \mathbf{X}$ tiene media y varianza dadas por, $\lambda = \mathbb{E}(Y) = \mathbf{t}^\top \boldsymbol{\mu}$ y $\sigma^2 = \text{var}(Y) = \mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t} \geq 0$ y tomando

$$\varphi_Y(\mathbf{t}) = \varphi_{t^\top \mathbf{X}}(1) = \exp(i\lambda - \frac{1}{2}\sigma^2) = \exp(i\mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t})$$

Como un caso particular tenemos que la función característica para $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \sigma^2 \mathbf{I}_p)$, es

$$\varphi_{\mathbf{Z}}(\mathbf{t}) = \exp(-\frac{1}{2}\sigma^2 \mathbf{t}^\top \mathbf{t}) = \prod_{i=1}^p \exp(-\frac{1}{2}\sigma^2 t_i^2) = \prod_{i=1}^p \varphi_{Z_i}(t_i)$$

y de este modo, se tiene que

$$\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \sigma^2 \mathbf{I}_p) \iff Z_1, \dots, Z_p \text{ IID } \mathbf{N}(0, \sigma^2).$$

RESULTADO 1.44. Suponga $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ y considere la transformación lineal

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b},$$

donde $\mathbf{A} \in \mathbb{R}^{m \times p}$ con $\text{rg}(\mathbf{A}) = m$. Entonces $\mathbf{Y} \sim \mathbf{N}_m(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$.

DEMOSTRACIÓN. Para mostrar el resultado usaremos la función característica de \mathbf{X} . Note que

$$\begin{aligned} \varphi_Y(\mathbf{t}) &= \mathbb{E}\{\exp(i\mathbf{t}^\top \mathbf{Y})\} = \mathbb{E}\{\exp(i\mathbf{t}^\top (\mathbf{A}\mathbf{X} + \mathbf{b}))\} \\ &= \exp(i\mathbf{t}^\top \mathbf{b}) \mathbb{E}\{\exp(i\mathbf{t}^\top \mathbf{A}\mathbf{X})\} \end{aligned}$$

Sea $\mathbf{h} = \mathbf{t}^\top \mathbf{A}$, entonces

$$\begin{aligned} \varphi_Y(\mathbf{t}) &= \exp(i\mathbf{t}^\top \mathbf{b}) \mathbb{E}\left(i\mathbf{h}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{h}^\top \boldsymbol{\Sigma} \mathbf{h}\right) = \exp(i\mathbf{t}^\top \mathbf{b}) \mathbb{E}\left(i\mathbf{t}^\top \mathbf{A}\boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top \mathbf{t}\right) \\ &= \exp\left(i\mathbf{t}^\top (\mathbf{A}\boldsymbol{\mu} + \mathbf{b}) - \frac{1}{2}\mathbf{t}^\top \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top \mathbf{t}\right), \end{aligned}$$

y el resultado sigue. □

RESULTADO 1.45. Si $\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ y $\boldsymbol{\Sigma}$ es definida positiva, entonces la densidad de \mathbf{X} es

$$f(\mathbf{x}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}.$$

DEMOSTRACIÓN. Sea Z_1, \dots, Z_p variables aleatorias IID $\mathcal{N}(0, 1)$. Entonces la densidad conjunta de $\mathbf{Z} = (Z_1, \dots, Z_p)^\top$ es

$$f(\mathbf{z}) = \prod_{i=1}^p (2\pi)^{-1/2} \exp(-z_i^2/2) = (2\pi)^{-p/2} \exp(-\frac{1}{2}\|\mathbf{z}\|^2).$$

Considere $\mathbf{X} = \boldsymbol{\mu} + \mathbf{B}\mathbf{Z}$ con $\boldsymbol{\mu} \in \mathbb{R}^p$ y $\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}^\top$, con \mathbf{B} matriz de rango completo. Entonces, tenemos la transformación inversa

$$\mathbf{Z} = \mathbf{g}^{-1}(\mathbf{X}) = \mathbf{B}^{-1}(\mathbf{X} - \boldsymbol{\mu}),$$

y $d\mathbf{Z} = d\mathbf{g}^{-1}(\mathbf{X}) = \mathbf{B}^{-1} d\mathbf{X}$, con matriz jacobiana $D\mathbf{g}^{-1}(\mathbf{X}) = \mathbf{B}^{-1}$, como

$$|D\mathbf{g}^{-1}(\mathbf{X})|_+ = |\mathbf{B}|^{-1} = |\mathbf{B}\mathbf{B}^\top|^{-1/2},$$

obtenemos

$$\begin{aligned} f(\mathbf{x}) &= |D\mathbf{g}^{-1}(\mathbf{x})|_+ f_{\mathbf{Z}}(\mathbf{g}^{-1}(\mathbf{x})) \\ &= (2\pi)^{-p/2} |\mathbf{B}\mathbf{B}^\top|^{-1/2} \exp\left\{\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{B}^{-\top} \mathbf{B}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}, \end{aligned}$$

notando que $\boldsymbol{\Sigma}^{-1} = \mathbf{B}^{-\top} \mathbf{B}^{-1}$ sigue el resultado deseado. \square

EJEMPLO 1.46. Sea $\mathbf{X} \sim \mathcal{N}_2(\mathbf{0}, \boldsymbol{\Sigma})$ donde

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}, \quad -1 < \rho < 1.$$

A continuación se presenta la función de densidad para los casos $\rho = 0.0, 0.4$ y 0.8 .

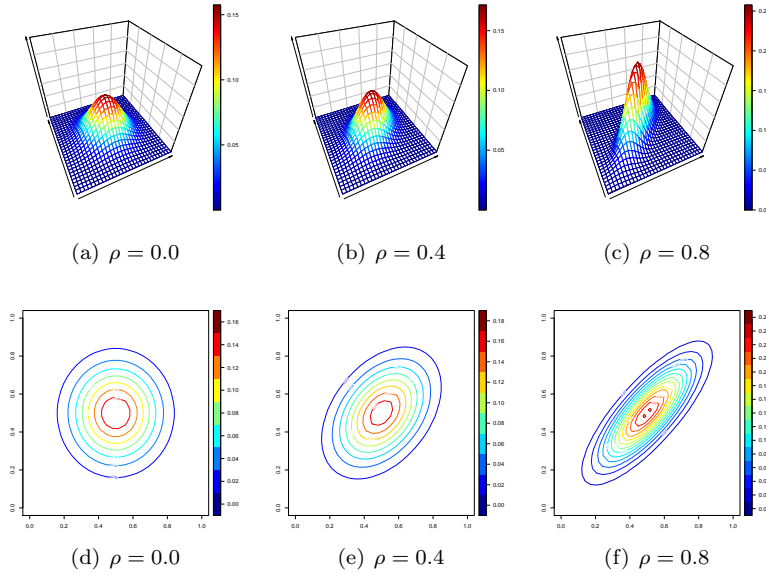


Figura 1. Densidad de $\mathbf{X} \sim \mathcal{N}_2(\mathbf{0}, \boldsymbol{\Sigma})$ para $\rho = 0.0, 0.4$ y 0.8 .

Note que la función de densidad es constante sobre el elipsoide en \mathbb{R}^p

$$(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = k,$$

para todo $k > 0$. Este elipsoide tiene centro $\boldsymbol{\mu}$, mientras que $\boldsymbol{\Sigma}$ determina su forma y orientación.

DEFINICIÓN 1.47. Sea \mathbf{U} vector aleatorio $p \times 1$ con distribución uniforme sobre el conjunto

$$\mathcal{S}_p = \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\| = 1\},$$

la superficie de la esfera unitaria en \mathbb{R}^p , y anotamos

$$\mathbf{U} \sim \mathcal{U}(\mathcal{S}_p).$$

La siguiente definición establece una relación entre la distribución normal y la uniforme sobre la esfera unitaria

PROPIEDAD 1.48. Si Z_1, \dots, Z_p son variables aleatorias IID con distribución $\mathcal{N}(0, 1)$, entonces $\mathbf{U} = (U_1, \dots, U_p)^\top$, definido como

$$\mathbf{U} = \frac{\mathbf{Z}}{\|\mathbf{Z}\|},$$

tiene distribución uniforme sobre la esfera unitaria \mathcal{S}_p .

El resultado anterior es muy relevante pues permite definir la densidad de un vector aleatorio $\mathbf{U} \sim \mathcal{U}(\mathcal{S}_p)$ y ofrece un procedimiento muy simple para generar observaciones sobre la esfera unitaria. Considere el siguiente gráfico,

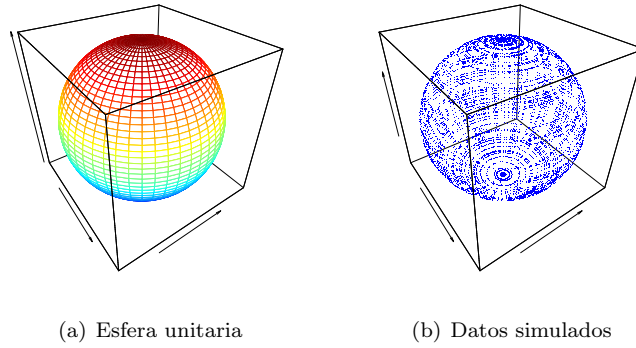


Figura 2. Esfera unitaria y datos generados sobre la superficie \mathcal{S}_p .

DEFINICIÓN 1.49. Un vector aleatorio $p \times 1$, \mathbf{X} se dice que tiene simetría esférica si para cualquier $\mathbf{Q} \in \mathcal{O}_p$, sigue que

$$\mathbf{Q}\mathbf{X} \stackrel{d}{=} \mathbf{X}.$$

EJEMPLO 1.50. Sea $\mathbf{U} \sim \mathcal{U}(\mathcal{S}_p)$, entonces es bastante obvio que $\mathbf{Q}\mathbf{U} \stackrel{d}{=} \mathbf{U}$.

DEFINICIÓN 1.51. Un vector aleatorio p -dimensional tiene distribución esférica sólo si su función característica satisface

$$(a) \quad \varphi(\mathbf{Q}^\top \mathbf{t}) = \varphi(\mathbf{t}), \text{ para todo } \mathbf{Q} \in \mathcal{O}_p.$$

(b) Existe una función $\phi(\cdot)$ de una variable escalar tal que $\varphi(\mathbf{t}) = \phi(\mathbf{t}^\top \mathbf{t})$. En este caso escribimos $\mathbf{X} \sim \mathcal{S}_p(\phi)$.

EJEMPLO 1.52. Sea $\mathbf{X} \sim \mathcal{N}_p(\mathbf{0}, \mathbf{I})$, tenemos que

$$\varphi(\mathbf{t}) = \exp\{-\frac{1}{2}(t_1^2 + \cdots + t_p^2)\} = \exp(-\frac{1}{2}\mathbf{t}^\top \mathbf{t}).$$

RESULTADO 1.53. Sea $\psi(\mathbf{t}^\top \mathbf{t})$ la función característica del vector aleatorio \mathbf{X} . Entonces \mathbf{X} tiene representación estocástica

$$\mathbf{X} \stackrel{d}{=} R\mathbf{U},$$

donde $\mathbf{U} \sim \mathcal{U}(\mathcal{S}_p)$ y $R \sim F(\mathbf{X})$ son independientes.

RESULTADO 1.54. Suponga que $\mathbf{X} \stackrel{d}{=} R\mathbf{U} \sim \mathcal{S}_p(\phi)$ ($P(\mathbf{X} = \mathbf{0}) = 0$), entonces

$$\|\mathbf{X}\| \stackrel{d}{=} R, \quad \frac{\mathbf{X}}{\|\mathbf{X}\|} \stackrel{d}{=} \mathbf{U}.$$

Además $\|\mathbf{X}\|$ y $\mathbf{X}/\|\mathbf{X}\|$ son independientes.

DEFINICIÓN 1.55 (Distribución de contornos elípticos). Un vector aleatorio $p \times 1$, \mathbf{X} tiene distribución de *contornos elípticos* con parámetros $\boldsymbol{\mu} \in \mathbb{R}^p$ y $\boldsymbol{\Sigma} \geq 0$, si

$$\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + \mathbf{B}\mathbf{Y}, \quad \mathbf{Y} \sim \mathcal{S}_k(\phi),$$

donde $\mathbf{B} \in \mathbb{R}^{p \times k}$ es matriz de rango completo tal que, $\mathbf{B}\mathbf{B}^\top = \boldsymbol{\Sigma}$ con $\text{rg}(\boldsymbol{\Sigma}) = k$. En cuyo caso escribimos $\mathbf{X} \sim \mathcal{EC}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \phi)$.

OBSERVACIÓN. La función característica de $\mathbf{X} \sim \mathcal{EC}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \phi)$ es de la forma:

$$\varphi(\mathbf{t}) = \exp(i\mathbf{t}^\top \boldsymbol{\mu})\phi(\mathbf{t}^\top \boldsymbol{\Sigma}\mathbf{t}).$$

Note además que la representación estocástica de \mathbf{X} es dada por

$$\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + R\mathbf{B}\mathbf{U},$$

donde $R \geq 0$ es independiente de \mathbf{U} y $\mathbf{B}\mathbf{B}^\top = \boldsymbol{\Sigma}$.

DEFINICIÓN 1.56. Se dice que el vector \mathbf{X} tiene distribución de contornos elípticos si su función de densidad es de la forma

$$f(\mathbf{x}) = |\boldsymbol{\Sigma}|^{-1/2} g((\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})), \quad \mathbf{x} \in \mathbb{R}^p,$$

donde $g: \mathbb{R} \rightarrow [0, \infty)$ es función decreciente, llamada *función generadora de densidad*, tal que:

$$\int_0^\infty u^{p/2-1} g(u) du < \infty.$$

OBSERVACIÓN. Asuma que $\mathbf{X} \sim \mathcal{EC}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \phi)$ con $\text{rg}(\boldsymbol{\Sigma}) = k$. Entonces,

$$Q(\mathbf{X}) = (\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^- (\mathbf{X} - \boldsymbol{\mu}) \stackrel{d}{=} R^2,$$

donde $\boldsymbol{\Sigma}^-$ es una inversa generalizada de $\boldsymbol{\Sigma}$.

EJEMPLO 1.57 (Distribución t multivariada). La función generadora de densidad de un vector aleatorio con distribución t multivariada asume la forma

$$g(u) = \frac{\Gamma(\frac{\nu+p}{2})}{\Gamma(\frac{\nu}{2})(\pi\nu)^{p/2}} \left(1 + \frac{u}{\nu}\right)^{-(\nu+p)/2}, \quad \nu > 0.$$

En cuyo caso escribimos, $\mathbf{X} \sim t_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$. Para esta distribución, tenemos que $R^2/p \sim F_{p, \nu}$.

EJEMPLO 1.58 (Distribución Exponencial Potencia). Para la distribución Exponencial Potencia (Gómez et al., 1988), la función generadora de densidades es dada por

$$g(u) = \frac{p\Gamma(\frac{p}{2})\pi^{-p/2}}{\Gamma(1 + \frac{p}{2\lambda})2^{1+\frac{p}{2\lambda}}} \exp(-u^\lambda/2), \quad \lambda > 0.$$

y es usual utilizar la notación $\mathbf{X} \sim \text{PE}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \lambda)$. En este caso tenemos que la variable aleatoria positiva R tiene densidad

$$h(r) = \frac{p}{\Gamma(1 + \frac{p}{2\lambda})2^{\frac{p}{2\lambda}}} r^{p-1} \exp(-r^{2\lambda}/2), \quad r > 0.$$

Note también que $R^{2\lambda} \sim \text{Gama}(\frac{1}{2}, \frac{p}{2\lambda})$.

EJEMPLO 1.59. En la siguiente figura se presenta la densidad asociadas a las siguientes funciones g :

- Normal: $g(u) = c_1 \exp(-u/2)$.
- Laplace: $g(u) = c_2 \exp(-\sqrt{u}/2)$.
- Cauchy: $g(u) = c_3(1+u)^{-(p+1)/2}$.
- Exponencial potencia (PE): $g(u) = c_4 \exp(-u^\lambda/2)$, $\lambda = 2$.

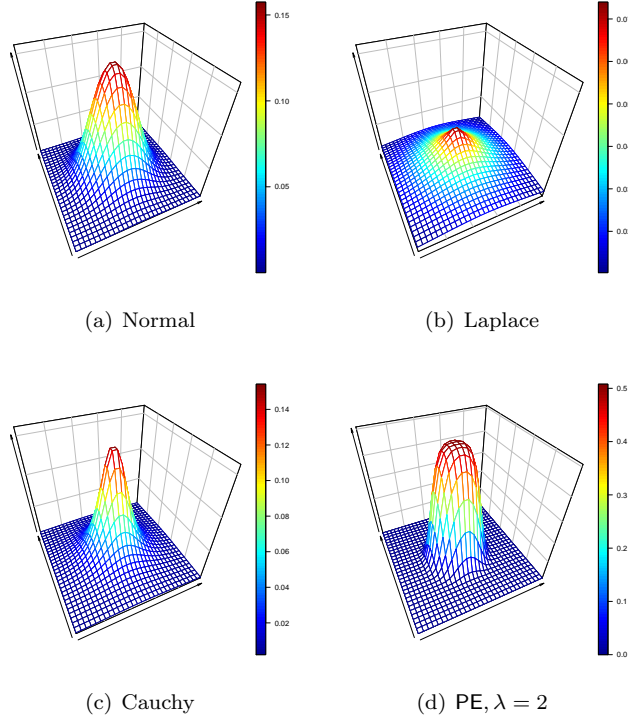


Figura 3. Funciones de densidad del vector $\mathbf{X} \sim \text{EC}_2(\mathbf{0}, \mathbf{I}; g)$ para las distribuciones normal, Laplace, Cauchy y exponencial potencia con $\lambda = 2$.

DEFINICIÓN 1.60 (Distribución de mezcla de escala normal). Sea $\boldsymbol{\mu} \in \mathbb{R}^p$, $\boldsymbol{\Sigma}$ matriz $p \times p$ definida positiva y H función de distribución de una variable aleatoria positiva, W . Entonces, se dice que el vector aleatorio \mathbf{X} sigue una *distribución de mezcla de escala normal* si su función de densidad asume la forma

$$f(\mathbf{x}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \int_0^\infty w^{p/2} \exp(-wu/2) dH(w),$$

donde $u = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$ y anotamos $\mathbf{X} \sim \text{SMN}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; H)$.

Esta familia de distribuciones ha sido frecuentemente sugerida como una interesante alternativa para producir estimadores robustos, manteniendo la elegancia y simplicidad de la teoría de máxima verosimilitud.

EJEMPLO 1.61 (Distribución Slash). Un vector aleatorio \mathbf{X} tiene distribución Slash si su función de densidad es de la forma:

$$f(\mathbf{x}) = \nu(2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \int_0^1 w^{p/2+\nu-1} \exp(-wu/2) dw.$$

En este caso, tenemos que $h(w) = \nu w^{\nu-1}$, para $w \in (0, 1)$ y $\nu > 0$. Es decir, $W \sim \text{Beta}(\nu, 1)$.

OBSERVACIÓN. Un vector aleatorio $\mathbf{X} \sim \text{SMN}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; H)$ si admite la representación

$$\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + W^{-1/2} \mathbf{Z},$$

donde $\mathbf{Z} \sim \text{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$ y $W \sim H(\delta)$ son independientes. Equivalentemente, podemos utilizar la estructura jerárquica:

$$\mathbf{X}|W \sim \text{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}/w), \quad W \sim H(\delta).$$

Por ejemplo, un vector aleatorio $\mathbf{X} \sim t_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$ puede ser caracterizado como:

$$\mathbf{X}|W \sim \text{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}/w), \quad W \sim \text{Gamma}(\nu/2, \nu/2).$$

Elementos de Inferencia

2.1. 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 2.1 (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 2.2. 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 2.3 (Factorización de Fisher-Neyman). *Suponga que X_1, \dots, X_n tiene densidad conjunta $f(\mathbf{x}; \theta)$, $\theta \in \Theta$. Una estadística $T : \mathcal{X}^n \rightarrow \mathbb{R}$ es suficiente para θ si y solo si, existe $g : \mathbb{R} \times \Theta \rightarrow \mathbb{R}$ y $h : \mathcal{X} \rightarrow \mathbb{R}$ tal que*

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

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 2.4. 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 2.5. 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) .

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

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

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

2.2. Función de verosimilitud

Para motivar ideas, suponga X_1, \dots, X_n variables aleatorias IID desde una función de distribución desconocida $G(x)$. Asumiremos que $G(x)$ corresponde al modelo estadístico verdadero (o verdadera CDF). Además, sea $F(x)$ un modelo arbitrario (el modelo asumido). Supondremos también que asociadas a $G(x)$ y $F(x)$ 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.

La información de Kullback-Leibler (KL) entre las funciones de densidad $g(x)$ y $f(x)$ es dada por:¹

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

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).

EJEMPLO 2.6. 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

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

mientras que (basta notar que $E_G[(X - \theta)^2] = \phi^2$):

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

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

$$I(g : f) = E_G(\log g(x)) - E_G(\log f(x)) = \frac{1}{2} \left\{ \log \frac{\sigma^2}{\phi^2} + \frac{\phi^2 + (\theta - \mu)^2}{\sigma^2} - 1 \right\}$$

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

Además, como

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

¹En ocasiones anotamos $I(G : F) = \int \log(g/f) dG$.

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

Claramente, mientras mayor sea este valor, mejor será el modelo (será más cercano a g). Además,

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

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

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

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

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

Esto lleva a la definición.

DEFINICIÓN 2.7 (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 2.8. Sea X_1, \dots, X_n variables aleatorias IID con distribución $\mathcal{N}(\theta, 1)$. 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 $\mathcal{N}_1(\bar{x}, 1/n)$.

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

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

de este modo, la variable aleatoria

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

Análogamente,

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

Es decir,

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

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

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

El siguiente lema permite justificar que aquél valor de $\boldsymbol{\theta}$ que se ajusta bien a los datos \mathbf{x} debe maximizar $L(\boldsymbol{\theta}; \mathbf{x})$.

LEMA 2.10 (Principio de verosimilitud). *Sea $\boldsymbol{\theta}_0$ el parámetro subyacente (o verdadero). Entonces,*

$$\mathbb{E}_{\boldsymbol{\theta}_0} \{\ell(\boldsymbol{\theta}; \mathbf{X})\} \leq \mathbb{E}_{\boldsymbol{\theta}_0} \{\ell(\boldsymbol{\theta}_0; \mathbf{X})\}, \quad \forall \boldsymbol{\theta}_0, \boldsymbol{\theta} \in \Theta.$$

DEMOSTRACIÓN. En efecto,²

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}_0} \{\log f(\mathbf{x}; \boldsymbol{\theta}_0)\} - \mathbb{E}_{\boldsymbol{\theta}_0} \{\log f(\mathbf{x}; \boldsymbol{\theta})\} &= -\mathbb{E}_{\boldsymbol{\theta}_0} \left\{ \log \frac{f(\mathbf{x}; \boldsymbol{\theta})}{f(\mathbf{x}; \boldsymbol{\theta}_0)} \right\} \\ &\geq -\log \mathbb{E}_{\boldsymbol{\theta}_0} \left\{ \frac{f(\mathbf{x}; \boldsymbol{\theta})}{f(\mathbf{x}; \boldsymbol{\theta}_0)} \right\} = \log \int \frac{f(\mathbf{x}; \boldsymbol{\theta})}{f(\mathbf{x}; \boldsymbol{\theta}_0)} f(\mathbf{x}; \boldsymbol{\theta}_0) d\mathbf{x} \\ &= \log \int f(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \log 1 = 0. \end{aligned}$$

□

²Recuerde que, para h convexa $\mathbb{E}(h(Y)) \geq h(\mathbb{E}(Y))$.

2.3. 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 $\{P_\theta : \theta \in \Theta\}$ tienen soporte común, de modo que el conjunto

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

no depende de θ .

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

EJEMPLO 2.11 (Contraejemplos). Considere $X \sim U(0, \theta)$, 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 $Y \sim \text{Exp}(a, b)$,

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

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 $\theta = (\theta_1, \dots, \theta_p)^\top$.

DEFINICIÓN 2.12 (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

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

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

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

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

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

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

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

RESULTADO 2.13. Bajo las condiciones A1 a A4, tenemos

$$E_\theta\{U(\theta; \mathbf{X})\} = \mathbf{0}, \quad \forall \theta \in \Theta$$

DEMOSTRACIÓN. Tenemos

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

□

EJEMPLO 2.14. Considere $X \sim N(\theta, 1)$. 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, $U(\theta; X) \sim N(0, 1)$. Así, es directo

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

DEFINICIÓN 2.15 (Matriz de información de Fisher). Suponga las condiciones **A1** a **A3**. Entonces la *matriz de información de Fisher* se define como:

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

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

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

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

EJEMPLO 2.16. Sean X_1, \dots, X_n variables aleatorias $N(\mu, \sigma^2)$ 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 2.17. 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

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

Por otro lado, note que

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

De este modo,

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

□

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,

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

DEFINICIÓN 2.18. 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 2.19. 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 $E(X) = n\theta$, sigue que

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

RESULTADO 2.20. Sean \mathbf{X} e \mathbf{Y} vectores aleatorios independientes con informaciones de Fisher $\mathcal{F}_X(\theta)$ y $\mathcal{F}_Y(\theta)$, respectivamente. Entonces la información de Fisher contenida en $\mathbf{Z} = (\mathbf{X}^{\top}, \mathbf{Y}^{\top})^{\top}$ es dada por

$$\mathcal{F}_Z(\theta) = \mathcal{F}_X(\theta) + \mathcal{F}_Y(\theta).$$

DEMOSTRACIÓN. Evidentemente, tenemos

$$\ell(\theta; \mathbf{z}) = \log f_1(\mathbf{x}; \theta) + \log f_2(\mathbf{y}; \theta),$$

mientras que

$$\mathbf{U}(\theta; \mathbf{z}) = \frac{\partial}{\partial \theta} \log f_1(\mathbf{x}; \theta) + \frac{\partial}{\partial \theta} \log f_2(\mathbf{y}; \theta) = \mathbf{U}(\theta; \mathbf{x}) + \mathbf{U}(\theta; \mathbf{y}),$$

por la independencia entre \mathbf{X} e \mathbf{Y} , sigue que

$$\begin{aligned}\mathcal{F}_Z(\theta) &= \text{Cov}(\mathbf{U}(\theta; \mathbf{X}) + \mathbf{U}(\theta; \mathbf{Y})) = \text{Cov}(\mathbf{U}(\theta; \mathbf{X})) + \text{Cov}(\mathbf{U}(\theta; \mathbf{Y})) \\ &= \mathcal{F}_X(\theta) + \mathcal{F}_Y(\theta).\end{aligned}$$

□

COROLARIO 2.21. Sea $\mathbf{X} = (X_1, \dots, X_n)^{\top}$ copias IID de una variable aleatoria X con información de Fisher (común) $\mathcal{F}_1(\theta)$. Entonces, la información contenida en la muestra es:

$$\mathcal{F}_n(\theta) = n\mathcal{F}_1(\theta).$$

OBSERVACIÓN. Suponga $\mathbf{X}_1, \dots, \mathbf{X}_n$ vectores aleatorios independientes (no necesariamente copias). Entonces,

$$\ell_n(\theta; \mathbf{X}) = \sum_{i=1}^n \ell_i(\theta; \mathbf{X}_i), \quad \mathbf{U}_n(\theta; \mathbf{X}) = \sum_{i=1}^n \mathbf{U}_i(\theta; \mathbf{X}_i), \quad \mathcal{F}_n(\theta) = \sum_{i=1}^n \mathcal{F}_i(\theta),$$

con $\mathbf{X} = (\mathbf{X}_1^{\top}, \dots, \mathbf{X}_n^{\top})^{\top}$.

RESULTADO 2.22 (Dependencia de la parametrización). Sea X una variable aleatoria con distribución P_{θ} . Suponga otra parametrización dada por $\theta = h(\phi)$ con h diferenciable. Entonces, la información contenida en X con relación a ϕ es dada por

$$\mathcal{F}_X^*(\phi) = \mathcal{F}_X(h(\phi))\{h'(\phi)\}^2,$$

donde $\mathcal{F}_X(\theta)$ es la información que $X \sim P_{\theta}$ contiene respecto de θ y $\mathcal{F}_X^*(\phi)$ es la información que $X \sim P_{h(\phi)}$ contiene sobre ϕ .

DEMOSTRACIÓN. La función de log-verosimilitud con relación a $P_{h(\phi)}$ es:

$$\ell^*(\phi; x) = \log f(x; h(\phi)).$$

De este modo, la función score asume la forma:

$$U^*(\phi; X) = \frac{\partial}{\partial \phi} \log f(x; h(\phi)) = \frac{\partial}{\partial \theta} \log f(x; \theta) \Big|_{\theta=h(\phi)} h'(\phi),$$

luego,

$$\begin{aligned}\mathcal{F}_X^*(\phi) &= \text{var}_\phi\{U^*(\phi; X)\} = \text{var}_\phi\{U(h(\phi); X) h'(\phi)\} \\ &= \{h'(\phi)\}^2 \text{var}_\phi\{U(h(\phi); X)\} = \{h'(\phi)\}^2 \mathcal{F}_X(h(\phi)).\end{aligned}$$

□

EJEMPLO 2.23. Considere la distribución $\text{Poi}(\theta)$. Así,

$$f(x; \theta) = \exp(x \log \theta - \theta) / x!$$

Note que el parámetro natural es $\eta = \log \theta$. Luego,

$$f(x; \eta) = \exp(x\eta - e^\eta) / x!$$

En la parametrización original, tenemos que

$$U(\theta; x) = \frac{\partial}{\partial \theta} \log f(x; \theta) = \frac{\partial}{\partial \theta} (x \log \theta - \theta) = \frac{x}{\theta} - 1.$$

De este modo,

$$\mathcal{F}_X(\theta) = \text{var}_\theta\{U(\theta; X)\} = \text{var}_\theta\left(\frac{X}{\theta} - 1\right) = \frac{1}{\theta^2} \text{var}_\theta(X) = \frac{1}{\theta}.$$

Por otro lado, $\theta = e^\eta = h(\eta)$. Así,

$$\mathcal{F}_X^*(\eta) = \mathcal{F}_X(h(\eta))\{h'(\eta)\}^2 = \frac{1}{e^\eta} (e^\eta)^2 = e^\eta.$$

OBSERVACIÓN. Considere $\boldsymbol{\theta} = \mathbf{h}(\boldsymbol{\phi})$ con $\mathbf{h} : \mathbb{R}^k \rightarrow \mathbb{R}^p$ ($p \leq k$) y sea $\mathbf{H}(\boldsymbol{\phi}) = \partial \mathbf{h}(\boldsymbol{\phi}) / \partial \boldsymbol{\phi}^\top$ matriz de rango completo. Entonces,

$$\mathcal{F}^*(\boldsymbol{\phi}) = \mathbf{H}^\top(\boldsymbol{\phi}) \mathcal{F}(\mathbf{h}(\boldsymbol{\phi})) \mathbf{H}(\boldsymbol{\phi}).$$

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 3.1. Una función $\mathbf{T} : \mathcal{X}^n \rightarrow \Gamma$ 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 $\mathbf{T}(\mathbf{x})$ es llamado estimación de $g(\theta)$ y corresponde a una realización de la variable aleatoria $\mathbf{T}(\mathbf{X})$.

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

3.1. Métodos de estimación

3.1.1. 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 3.2 (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

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

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 3.3. 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 3.4. 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 3.5. 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 3.6. 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 3.7. 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 3.8. Suponga X_1, \dots, X_n una muestra aleatoria desde $\text{Poi}(\lambda)$, $\lambda > 0$. Recuerde que

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

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

$$\begin{aligned} \theta_1 &= h_1(\mu_1, \dots, \mu_p), \\ &\vdots \\ \theta_p &= h_p(\mu_1, \dots, \mu_p). \end{aligned} \tag{3.1}$$

¹Más adelante estudiaremos como comparar entre dos estimadores.

Finalmente, el estimador $\hat{\boldsymbol{\theta}}_{\text{MM}}$, puede ser obtenido substituyendo en (3.1) 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:

$$\boldsymbol{\Psi}(\boldsymbol{\theta}) = \mathbf{0}, \quad (3.2)$$

donde $\boldsymbol{\Psi} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ y por tanto $\hat{\boldsymbol{\theta}}_{\text{MM}}$ corresponde a una raíz de la **ecuación de estimación** en (3.2). Note que, en general, para obtener $\hat{\boldsymbol{\theta}}_{\text{MM}}$ se requiere de métodos iterativos, tal como:

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

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

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

3.1.2. Estimación máximo verosímil. Este es uno de los procedimientos de estimación más ampliamente usados. Es motivado por el principio de verosimilitud (Lema 2.10) y los estimadores obtenidos disfrutan de buenas propiedades.

DEFINICIÓN 3.9 (Estimador máximo verosímil). Un estimador $\hat{\boldsymbol{\theta}}_{\text{ML}}$ es llamado estimador máximo verosímil (MLE) de $\boldsymbol{\theta}$, si

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

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

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

o equivalentemente,

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

Además, en ocasiones escribimos

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

RESULTADO 3.10 (Invarianza del MLE). Si $\boldsymbol{\gamma} = \mathbf{g}(\boldsymbol{\theta})$ y \mathbf{g} es biyectiva. Entonces $\hat{\boldsymbol{\theta}}$ es el MLE para $\boldsymbol{\theta}$ si y solo si $\hat{\boldsymbol{\gamma}} = \mathbf{g}(\hat{\boldsymbol{\theta}})$ es el MLE para $\boldsymbol{\gamma}$.

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 3.11. 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 3.12 (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 3.13 (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 3.14 (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 3.15 (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 3.16 (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 3.17 (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.

3.1.3. Implementación del método de máxima verosimilitud. Anteriormente hemos revisado ejemplos en los que ha sido posible obtener los estimadores ML de forma explícita. Sin embargo, en general es necesario recurrir a métodos numéricos, los que involucran la elección de una estimación inicial $\boldsymbol{\theta}^{(0)}$ y sucesivamente se construye la secuencia $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots$, de manera de alcanzar convergencia a la solución $\hat{\boldsymbol{\theta}}_{\text{ML}}$.

Con el objetivo de resolver el problema

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

suponga la expansión de Taylor en torno de $\boldsymbol{\theta}^*$, como:

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

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,

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

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

$$\mathbf{H}(\boldsymbol{\theta}^{(k)}) \mathbf{p} = -\mathbf{U}(\boldsymbol{\theta}^{(k)}). \quad (3.3)$$

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 (3.3), 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 (3.3), 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. De esta manera el procedimiento para determinar una aproximación de la solución asume la forma:

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + \lambda_k \{ -\mathbf{H}(\boldsymbol{\theta}^{(k)}) \}^{-1} \mathbf{U}(\boldsymbol{\theta}^{(k)}), \quad k = 0, 1, \dots,$$

que es conocido como *método de Newton-Raphson*. Es conocido que este procedimiento converge rápidamente siempre que un valor inicial cercano al óptimo haya sido escogido. En casos en que sea difícil llevar a cabo el cálculo de la matriz Hessiana podemos utilizar un *método quasi-Newton*, el cual no involucra información de segundo orden sino que utiliza la dirección de búsqueda:

$$\mathbf{p}_k = \{ -\mathbf{B}_k \}^{-1} \mathbf{U}(\boldsymbol{\theta}^{(k)}),$$

donde \mathbf{B}_k es una aproximación de la matriz Hessiana $\mathbf{H}(\boldsymbol{\theta}^{(k)})$. 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 algoritmo 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 $\mathbf{s}_k = \boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}$. Es frecuente usar como valor inicial para \mathbf{B}_0^{-1} la matriz identidad. En situaciones en que también es difícil calcular el vector score $\mathbf{U}(\boldsymbol{\theta})$ es posible determinar la información de primer orden usando diferenciación numérica.

Un procedimiento popular en estadística corresponde al *método Fisher-scoring*, que corresponde a un algoritmo en la clase quasi-Newton donde $-\mathbf{H}(\boldsymbol{\theta})$ es aproximada mediante la matriz de información de Fisher. De este modo, obtenemos el siguiente esquema iterativo:

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

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

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, $\mathbf{U}_n(\boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{U}_i(\boldsymbol{\theta})$, con $\mathbf{U}_i(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}_i; \boldsymbol{\theta})$. Notando que

$$\frac{1}{n} \sum_{i=1}^n \mathbf{U}_i(\boldsymbol{\theta}) \mathbf{U}_i^\top(\boldsymbol{\theta}) \xrightarrow{P} \mathbb{E}\{\mathbf{U}_n(\boldsymbol{\theta}) \mathbf{U}_n^\top(\boldsymbol{\theta})\},$$

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

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

Usualmente, llevamos a cabo la iteración del procedimiento de estimación hasta que

$$\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}\| < \tau,$$

donde τ , conocido como tolerancia, sea pequeño (por ejemplo, $\tau = 10^{-6}$). Otra alternativa es considerar el criterio

$$(\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)})^\top \mathcal{F}(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}) < \delta,$$

donde δ representa un valor de tolerancia. Por ejemplo, basados en la idea de un *elipsoide de confianza*, Bates y Watts (1981) propusieron un criterio de convergencia basado en un *offset relativo*. Tal idea puede ser aplicada a diversos modelos estadísticos. En efecto, ha sido implementada para los modelos disponibles en las bibliotecas para R nlme (Pinheiro et al., 2019) y heavy (Osorio, 2019).

EJEMPLO 3.18 (distribución Cauchy). Suponga X_1, \dots, X_n variables aleatorias desde la distribución $\text{Cauchy}(\theta, 1)$, con densidad

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

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

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

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

$$U(\theta; \mathbf{x}) = \sum_{i=1}^n \frac{2}{1 + (x_i - \theta)^2} (x_i - \theta) = \sum_{i=1}^n \omega_i(\theta)(x_i - \theta) = 0, \quad (3.4)$$

donde $\omega_i(\theta) = 2/(1 + (x_i - \theta)^2)$. Es fácil notar que la Ecuación (3.4) 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:

$$\begin{aligned} \theta^{(k+1)} &= \theta^{(k)} - \frac{U(\theta^{(k)}; \mathbf{x})}{H(\theta^{(k)}; \mathbf{x})} \\ &= \theta^{(k)} + \frac{\sum_{i=1}^n \omega_i(\theta^{(k)})(x_i - \theta^{(k)})}{\sum_{i=1}^n \omega_i(\theta^{(k)}) \{1 - \omega_i(\theta^{(k)})(x_i - \theta^{(k)})^2\}}. \end{aligned} \quad (3.5)$$

OBSERVACIÓN. En la biblioteca `heavy` se encuentra una alternativa al esquema iterativo en (3.5) que utiliza un *Algoritmo EM*.³

³El Algoritmo EM (o de Esperanza-Maximización) permite obtener los estimadores ML en presencia de *datos perdidos*.

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 3.19 (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 $\hat{\lambda}_{\text{ML}} = e^{\hat{\phi}_{\text{ML}}}$.

3.2. 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 3.20 (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 3.21 (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 3.22 (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 3.23. 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,⁴

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

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

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

⁴Para \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 γ es k -dimensional, esto es, $\mathbf{g} : \Theta \rightarrow \Gamma \subseteq \mathbb{R}^m$. Entonces, \mathbf{T} se dice un estimador insesgado, si

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

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

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

DEFINICIÓN 3.24. Sean \mathbf{T} y \mathbf{T}_* dos estimadores para γ . Decimos que \mathbf{T}_* tiene error cuadrático medio más pequeño que \mathbf{T} si

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

y escribimos

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

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 \mathbf{T}_* es *T-óptimo*, si

$$\text{tr MSE}(\mathbf{T}_*, \boldsymbol{\theta}) \leq \text{tr MSE}(\mathbf{T}, \boldsymbol{\theta}), \quad (3.6)$$

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

$$\det \text{MSE}(\mathbf{T}_*, \boldsymbol{\theta}) \leq \det \text{MSE}(\mathbf{T}, \boldsymbol{\theta}). \quad (3.7)$$

El criterio dado en la Definición 3.24 también es conocido como *M-optimalidad*. Debemos destacar que en ocasiones el error cuadrático medio es definido como

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

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.

Una alternativa es restringirse a alguna subclase de estimadores. De este modo, nos concentraremos en la clase de estimadores insesgados.

DEFINICIÓN 3.25 (Mejor estimador insesgado). Para el modelo estadístico $\mathcal{P} = \{\mathbf{P}_\theta : \theta \in \Theta\}$, un estimador insesgado T_* para $\gamma = g(\theta) \in \mathbb{R}$ se dice el mejor estimador insesgado (BUE), si para cualquier otro estimador insesgado

$$\text{var}_\theta(T_*) \leq \text{var}_\theta(T), \quad \forall \theta \in \Theta.$$

RESULTADO 3.26 (Cota de Cramér-Rao). *Suponga que se satisfacen las condiciones A1-A4, que la información de Fisher es tal que $0 < \mathcal{F}_X(\theta) < \infty$ y sea $\gamma = g(\theta)$, donde g es continua y diferenciable con $g' \neq 0$. Si T es estimador insesgado para γ , entonces*

$$\text{var}_\theta(T) \geq \frac{\{g'(\theta)\}^2}{\mathcal{F}_X(\theta)}, \quad \forall \theta \in \Theta.$$

DEMOSTRACIÓN. Considere

$$\begin{aligned} \text{Cov}_\theta(T(\mathbf{X}), U(\theta; \mathbf{X})) &= \mathbb{E}_\theta(T(\mathbf{X}) U(\theta; \mathbf{X})) = \int_A T(\mathbf{x}) \frac{\partial}{\partial \theta} \log f(\mathbf{x}; \theta) f(\mathbf{x}; \theta) d\mathbf{x} \\ &= \int_A T(\mathbf{x}) \frac{\partial}{\partial \theta} f(\mathbf{x}; \theta) d\mathbf{x}. \end{aligned}$$

Como T es regular⁵ y es insesgada, sigue que

$$\text{Cov}_\theta(T(\mathbf{X}), U(\theta; \mathbf{X})) = \frac{\partial}{\partial \theta} \int_A T(\mathbf{x}) f(\mathbf{x}; \theta) d\mathbf{x} = \frac{\partial}{\partial \theta} \mathbb{E}_\theta(T(\mathbf{X})) = g'(\theta).$$

Haciendo $h(\theta) = g'(\theta)/\mathcal{F}_X(\theta)$, tenemos

$$\begin{aligned} 0 &\leq \text{var}_\theta(T(\mathbf{X}) - h(\theta)U(\theta; \mathbf{X})) \\ &= \text{var}_\theta(T(\mathbf{X})) + h^2(\theta) \text{var}_\theta(U(\theta; \mathbf{X})) - 2h(\theta) \text{Cov}_\theta(T(\mathbf{X}), U(\theta; \mathbf{X})) \\ &= \text{var}_\theta(T(\mathbf{X})) + h^2(\theta) \mathcal{F}_X(\theta) - 2h(\theta)g'(\theta), \end{aligned}$$

es decir,

$$0 \leq \text{var}_\theta(T(\mathbf{X})) - \{g'(\theta)\}^2/\mathcal{F}_X(\theta).$$

□

DEFINICIÓN 3.27 (Eficiencia). La eficiencia de un estimador insesgado T es definida como la razón de su varianza y la cota de Cramér-Rao. Esto es,

$$\text{EFF}(T, \theta) = \frac{\{g'(\theta)\}^2/\mathcal{F}_X(\theta)}{\text{var}_\theta(T)}.$$

Un estimador que alcanza la cota de Cramér-Rao se dice un *estimador eficiente*. Más aún, un estimador eficiente es BUE.

EJEMPLO 3.28. Considere X_1, \dots, X_n variables aleatorias IID desde $\text{Exp}(\lambda)$ con $\lambda > 0$. Sea $\gamma = 1/\lambda$ el parámetro de interés. Un estimador insesgado para λ es \bar{X} con varianza $\frac{1}{n\lambda^2}$. Como la información de Fisher es n/λ^2 y $g'(\lambda) = -1/\lambda^2$, tenemos

$$\frac{\{-1/\lambda^2\}^2}{n/\lambda^2} = \frac{1}{n\lambda^2},$$

es decir, \bar{X} es eficiente.

⁵Es decir, podemos intercambiar las operaciones de integración y diferenciación.

Considere el caso en que $\Theta \subset \mathbb{R}^k$ y que el parámetro de interés es $\gamma \in \mathbb{R}^m$ con $\gamma = g(\theta)$. Sea \mathbf{T} y \mathbf{T}_* dos estimadores insesgados de γ . Decimos que \mathbf{T}_* tiene covarianza más pequeña que \mathbf{T} , si

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

en cuyo caso escribimos

$$\text{Cov}_\theta(\mathbf{T}_*) \leq \text{Cov}_\theta(\mathbf{T}). \quad (3.8)$$

Evidentemente, Ecuación (3.8) corresponde a un caso particular de la Definición 3.24 para el caso de estimadores insesgados. Esto permite extender la cota de Cramér-Rao para el caso multiparamétrico.

Suponga que las condiciones A1 a A4 son satisfechas y que la matriz de información de Fisher es no singular. Entonces la cota de Cramér-Rao asume la forma:

$$\text{Cov}_\theta(\mathbf{T}) \geq \left(\frac{\partial g(\theta)}{\partial \theta^\top} \right) \mathcal{F}_X^{-1}(\theta) \left(\frac{\partial g(\theta)}{\partial \theta^\top} \right)^\top, \quad \forall \theta \in \Theta.$$

EJEMPLO 3.29. Suponga que tenemos X_1, \dots, X_n variables aleatorias independientes desde $N(\mu, \sigma^2)$ con parámetro de interés $\theta = (\mu, \sigma^2)^\top$. La matrix de información de Fisher está dada por

$$\mathcal{F}_X(\theta) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix},$$

cuya matriz inversa corresponde a la cota de Cramér-Rao. Sabemos que los estimadores insesgados \bar{X} y S^2 son independientes y que

$$\text{var}_\theta(\bar{X}) = \frac{\sigma^2}{n}, \quad \text{var}_\theta(S^2) = \frac{2\sigma^4}{n-1}.$$

Por tanto no se alcanza la cota inferior para la varianza.

RESULTADO 3.30. *Bajo las condiciones A1 a A4 el estimador de máxima verosimilitud satisface las siguientes propiedades:*

- i) *El estimador ML depende de los datos via la estadística suficiente.*
- ii) *Si existe un estimador insesgado y eficiente $\tilde{\theta}$, entonces $\tilde{\theta} = \hat{\theta}_{\text{ML}}$.*

DEMOSTRACIÓN. i) sigue notando que, por el Teorema de factorización de Fisher-Neyman

$$L(\theta; \mathbf{x}) \propto g(\mathbf{T}(\mathbf{x}), \theta),$$

para \mathbf{T} suficiente.

Por simplicidad para la prueba de ii) sólo consideraremos el caso en que θ es real-valuado. Sabemos que $\tilde{\theta}$ alcanza la cota de Cramér-Rao pues es un estimador eficiente. Así, tenemos que

$$\tilde{\theta}(\mathbf{x}) - \theta = \frac{U(\theta; \mathbf{x})}{\mathcal{F}_X(\theta)}, \quad \forall \theta \in \Theta,$$

que es válido en particular para $\theta = \hat{\theta}_{\text{ML}}$. Es decir,

$$\tilde{\theta}(\mathbf{x}) - \hat{\theta}_{\text{ML}} = \frac{U(\hat{\theta}_{\text{ML}}; \mathbf{x})}{\mathcal{F}_X(\hat{\theta}_{\text{ML}})}.$$

Como $\hat{\theta}_{\text{ML}}$ maximiza $\ell(\theta; \mathbf{x})$ y por tanto, $U(\hat{\theta}_{\text{ML}}; \mathbf{x}) = 0$. De este modo,

$$\tilde{\theta}(\mathbf{x}) - \hat{\theta}_{\text{ML}} = 0.$$

□

3.3. Propiedades Asintóticas

Parece razonable que conforme el tamaño muestral crece mayor confianza tendremos en nuestras inferencias, debido a que la muestra contendrá más información con respecto a la distribución subyacente. Para caracterizar las propiedades asintóticas de los estimadores se introducirá algunas nociones de convergencia de variables aleatorias.

DEFINICIÓN 3.31 (Consistencia). Decimos que una secuencia de estimadores $\{T_n\}$ para el parámetro $\gamma = g(\theta)$ es *débilmente consistente*, si T_n converge en probabilidad a γ , esto es, si para cualquier $\epsilon > 0$ y para todo $\theta \in \Theta$,

$$\lim_{n \rightarrow \infty} P_\theta(|T_n - g(\theta)| > \epsilon) = 0,$$

y escribimos $T_n \xrightarrow{P} \gamma$. Mientras que, se dice que $\{T_n\}$ converge con probabilidad 1 o casi seguramente (a.s.) a γ , para todo $\theta \in \Theta$, esto es,

$$P_\theta \left(\lim_{n \rightarrow \infty} T_n = g(\theta) \right) = 1,$$

es decir T_n es *consistente fuerte*, en cuyo caso anotamos $T_n \xrightarrow{\text{a.s.}} \gamma$.

Consistencia muy frecuentemente es una consecuencia de la ley de los grandes números. Es tipo más simple es la convergencia de la media muestral. Por el teorema de mapeo continuo es posible obtener la consistencia de otros estimadores.

TEOREMA 3.32 (Teorema del mapeo continuo). Sea $\{S_n\}$ una secuencia de variables aleatorias, S_0 una variable aleatoria y h una función continua.

i) Si $S_n \xrightarrow{P} S_0$, entonces

$$h(S_n) \xrightarrow{P} h(S_0).$$

ii) Si $S_n \xrightarrow{\text{a.s.}} S_0$, entonces

$$h(S_n) \xrightarrow{\text{a.s.}} h(S_0).$$

Una herramienta importante para la verificación de consistencia es la desigualdad de Chebyshev. Para una variable aleatoria Z con media finita, tenemos

$$P(|Z - E(Z)| > \tau) \leq \frac{\text{var}(Z)}{\tau^2}.$$

EJEMPLO 3.33. Sea $\{X_n\}$ una secuencia de variables aleatorias IID con función de distribución F . Entonces se tiene:

1. La media aritmética converge a $E_F(X) = \mu$, es decir, $\bar{X}_n \xrightarrow{P} \mu$. En efecto, usando la desigualdad de Chebyshev y asumiendo $\sigma < \infty$, tenemos

$$P(|\bar{X}_n - \mu| > \epsilon) \leq \frac{\text{var}(\bar{X}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \rightarrow 0,$$

para $n \rightarrow \infty$.

2. La varianza empírica y la desviación estándar convergen a $\text{var}_F(X) = \sigma^2$ y σ , respectivamente:

$$S_n^2 \xrightarrow{P} \sigma^2, \quad S_n \xrightarrow{P} \sigma.$$

3. La frecuencia relativa de un evento A converge a su probabilidad. Sea

$$Q_n(A) = \frac{1}{n} \sum_{i=1}^n I(X_i \in A),$$

donde $I(\cdot)$ denota la función indicadora. Entonces $Q_n(A) \xrightarrow{P} P(A)$.

4. Si el j -ésimo momento $\mu_j = E_F(X^j)$ existe, entonces los momentos empíricos

$$m_j = \frac{1}{n} \sum_{i=1}^n X_i^j,$$

son consistentes. Es decir, $m_j \xrightarrow{P} \mu_j$. De ahí que, el estimador de momentos para $\gamma = h(\mu_1, \dots, \mu_r)$ es un estimador consistente si la función h es continua.

DEFINICIÓN 3.34 (Convergencia en distribución). Sea $\{X_n\}$ una secuencia de variables aleatorias y sea X otra variable aleatoria. Además, considere F_n la CDF de X_n y F la CDF de X . Se dice que X_n converge en distribución a X , en cuyo caso escribimos $X_n \xrightarrow{D} X$ si

$$\lim_{n \rightarrow \infty} F_n(t) = F(t),$$

para todo t donde F es continua.

TEOREMA 3.35. Sea $\{Z_n\}$ una secuencia de variables aleatorias y sea $M_n(t)$ la MGF de Z_n . Sea Z una variable aleatoria con MGF dada por $M(t)$. Si

$$M_n(t) \rightarrow M(t), \quad \text{para todo } |t| < h, h > 0.$$

Entonces, $Z_n \xrightarrow{D} Z$.

TEOREMA 3.36 (Teorema de Slutsky). Considere dos secuencias de variables aleatorias $\{X_n\}$, $\{Y_n\}$, una variable aleatoria X y una constante fija c . Suponga que $X_n \xrightarrow{D} X$, y $Y_n \xrightarrow{P} c$. Entonces:

$$i) X_n \pm Y_n \xrightarrow{D} X \pm c.$$

$$ii) X_n Y_n \xrightarrow{D} cX.$$

$$iii) X_n Y_n^{-1} \xrightarrow{D} c^{-1}X \text{ siempre que } P(Y_n = 0) = 0 \text{ para todo } n \text{ y } c \neq 0.$$

DEFINICIÓN 3.37 (Normalidad asintótica). Una secuencia de estimadores $\{\mathbf{T}_n\}$ para el parámetro m -dimensional $\gamma = \mathbf{g}(\boldsymbol{\theta})$ es *asintóticamente normal* si para todo $\boldsymbol{\theta} \in \Theta$ la distribución de $\sqrt{n}(\mathbf{T}_n - \mathbf{g}(\boldsymbol{\theta}))$ converge a una distribución normal con media cero y matriz de covarianza $\boldsymbol{\Sigma}(\boldsymbol{\theta})$. Es decir,

$$\sqrt{n}(\mathbf{T}_n - \mathbf{g}(\boldsymbol{\theta})) \xrightarrow{D} N_m(\mathbf{0}, \boldsymbol{\Sigma}(\boldsymbol{\theta})).$$

La normalidad asintótica permite establecer si un estimador es *asintóticamente eficiente*. En efecto, diremos que un estimador es asintóticamente eficiente si este es asintóticamente normal, con

$$\boldsymbol{\Sigma}(\boldsymbol{\theta}) = \left(\frac{\partial \mathbf{g}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right) \mathcal{F}_X^{-1}(\boldsymbol{\theta}) \left(\frac{\partial \mathbf{g}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right)^\top,$$

donde $\mathcal{F}_X(\boldsymbol{\theta})$ es la matriz de información de la distribución subyacente.

RESULTADO 3.38 (Método Delta). Suponga T_n un estimador de la forma $T_n = h(S_n)$ donde la secuencia $\{S_n\}$ es asintóticamente normal, esto es,

$$\sqrt{n}(S_n - \mu) \xrightarrow{D} N(\mathbf{0}, \Sigma),$$

para $\mu \in \mathbb{R}^k$ y $\Sigma > 0$. Si $\partial h(\mu)/\partial \mu^\top$ es matriz de rango completo, entonces

$$\sqrt{n}(T_n - h(\mu)) \xrightarrow{D} N\left(\mathbf{0}, \left(\frac{\partial h(\mu)}{\partial \mu^\top}\right) \Sigma \left(\frac{\partial h(\mu)}{\partial \mu^\top}\right)^\top\right).$$

EJEMPLO 3.39. Sean X_1, \dots, X_n variables aleatorias IID con $E(X_i) = \mu \neq 0$ y $\text{var}(X_i) = \sigma^2 < \infty$. El parámetro $\gamma = \log \mu$ es estimado por $\hat{\gamma}_n = \log \bar{X}_n$. Este estimador es consistente y asintóticamente normal. En efecto, $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} N(0, \sigma^2)$. Como $h(s) = \log s$ y $h'(s) = 1/s$, sigue que

$$\sqrt{n}(\log \bar{X}_n - \log \mu) \xrightarrow{D} N\left(0, \frac{\sigma^2}{\mu^2}\right).$$

3.3.1. Distribución asintótica de los estimadores ML. Considere una secuencia de estimadores $\{\hat{\theta}_n\}$ que converge en probabilidad a θ_0 perteneciendo al interior de Θ , y suponga que

SUPUESTO A5. La matriz

$$\mathcal{F}_1(\theta_0) = E_{\theta_0} \left\{ -\frac{\partial^2 \log f(\mathbf{y}_1; \theta_0)}{\partial \theta \partial \theta^\top} \right\},$$

existe y es no singular.

RESULTADO 3.40. Bajo las condiciones **A1-A5** una secuencia de máximos locales de $\ell_n(\theta)$ tiene distribución asintótica

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N_k(\mathbf{0}, \mathcal{F}_1^{-1}(\theta_0)).$$

DEMOSTRACIÓN. Note que la secuencia $\{\hat{\theta}_n\}$ satisface las ecuaciones de verosimilitud, $\partial \ell_n(\theta)/\partial \theta = \mathbf{0}$. Usando una expansión de Taylor del vector score en torno de $\theta = \theta_0$, resulta⁶

$$\frac{\partial \ell_n(\theta)}{\partial \theta} = \frac{\partial \ell_n(\theta_0)}{\partial \theta} + \frac{\partial^2 \ell_n(\theta_0)}{\partial \theta \partial \theta^\top} (\theta - \theta_0) + o_P(\mathbf{1}),$$

haciendo $\theta = \hat{\theta}_n$, sigue que

$$\mathbf{0} = \frac{\partial \ell_n(\theta_0)}{\partial \theta} + \frac{\partial^2 \ell_n(\theta_0)}{\partial \theta \partial \theta^\top} (\hat{\theta}_n - \theta_0) + o_P(\mathbf{1}),$$

es decir

$$-\frac{\partial^2 \ell_n(\theta_0)}{\partial \theta \partial \theta^\top} (\hat{\theta}_n - \theta_0) = \frac{\partial \ell_n(\theta_0)}{\partial \theta} + o_P(\mathbf{1}),$$

o equivalentemente,

$$\left(-\frac{1}{n} \frac{\partial^2 \ell_n(\theta_0)}{\partial \theta \partial \theta^\top} \right) \sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{1}{\sqrt{n}} \frac{\partial \ell_n(\theta_0)}{\partial \theta} + o_P(\mathbf{1}).$$

⁶Si $Z_n \xrightarrow{P} 0$ entonces anotamos $Z_n = o_P(1)$.

Por otro lado, tenemos que

$$-\frac{1}{n} \frac{\partial^2 \ell_n(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = \frac{1}{n} \sum_{i=1}^n -\frac{\partial^2 \log f(\mathbf{Y}_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$$

converge casi seguramente a

$$\mathcal{F}_1(\boldsymbol{\theta}_0) = \mathbb{E}_{\boldsymbol{\theta}_0} \left(-\frac{\partial^2 \log f(\mathbf{Y}_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right),$$

debido a la ley fuerte de los grandes números. Lo anterior lleva a,

$$\mathcal{F}_1(\boldsymbol{\theta}_0) \sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = \frac{1}{\sqrt{n}} \frac{\partial \ell_n(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} + o_P(\mathbf{1}),$$

Note también que,

$$\begin{aligned} \frac{1}{\sqrt{n}} \frac{\partial \ell_n(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial \log f(\mathbf{Y}_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \frac{\partial \log f(\mathbf{Y}_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} - \mathbb{E}_{\boldsymbol{\theta}_0} \left(\frac{\partial \log f(\mathbf{Y}_i; \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right) \right\}, \end{aligned}$$

converge en distribución a

$$\mathbf{N}_k(\mathbf{0}, \text{Cov}_{\boldsymbol{\theta}_0}(\mathbf{U}_i(\boldsymbol{\theta}_0))) \stackrel{d}{=} \mathbf{N}_k(\mathbf{0}, \mathcal{F}_1(\boldsymbol{\theta}_0)).$$

Premultiplicando por $\mathcal{F}_1^{-1}(\boldsymbol{\theta}_0)$ sigue que

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = \frac{1}{\sqrt{n}} \mathcal{F}_1^{-1}(\boldsymbol{\theta}_0) \mathbf{U}_n(\boldsymbol{\theta}_0) + o_P(\mathbf{1}).$$

Usando que

$$\frac{1}{\sqrt{n}} \mathcal{F}_1^{-1}(\boldsymbol{\theta}_0) \mathbf{U}_n(\boldsymbol{\theta}_0) \xrightarrow{D} \mathbf{N}_k(\mathbf{0}, \mathcal{F}_1^{-1}(\boldsymbol{\theta}_0) \mathcal{F}_1(\boldsymbol{\theta}_0) \mathcal{F}_1^{-1}(\boldsymbol{\theta}_0)),$$

sigue el resultado deseado. \square

El resultado anterior es de gran importancia, pues bajo condiciones de regularidad el MLE es consistente, asintóticamente normal y eficiente, es decir es BUE. Además, esto permite el desarrollo de intervalos de confianza y test de hipótesis asintóticos.

Intervalos y Regiones de Confianza

El objetivo de esta sección es abordar el problema $\theta \in C$, donde $C \subseteq \Theta$, $C = C(\mathbf{X})$ es un conjunto determinado por los datos observados $\mathbf{X} = \mathbf{x}$.

OBSERVACIÓN. Si θ es real-valuado, entonces C corresponde a un intervalo.

DEFINICIÓN 4.1. Una estimación intervalar de un parámetro real-valuado θ es cualquier par de funciones $L(x_1, \dots, x_n)$ y $U(x_1, \dots, x_n)$ que satisfacen

$$L(\mathbf{x}) \leq U(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}.$$

Para $\mathbf{X} = \mathbf{x}$ tenemos $L(\mathbf{x}) \leq \theta \leq U(\mathbf{x})$, mientras que $[L(\mathbf{X}), U(\mathbf{X})]$ es un intervalo aleatorio.

Aunque estamos interesados en intervalos de la forma $[L(\mathbf{x}), U(\mathbf{x})]$ también podemos tener $(-\infty, U(\mathbf{x})]$ o bien $[L(\mathbf{x}), \infty)$ para indicar $\theta \leq U(\mathbf{x})$ o $L(\mathbf{x}) \leq \theta$, respectivamente.

EJEMPLO 4.2. Considere X_1, X_2, X_3, X_4 una muestra aleatoria desde $N(\mu, 1)$. Un estimador intervalar de μ es $[\bar{X} - 1, \bar{X} + 1]$, es decir

$$\mu \in [\bar{X} - 1, \bar{X} + 1]$$

Note que $\bar{X} \sim N(\mu, 1/4)$, pero

$$P(\bar{X} = \mu) = 0$$

De ahí que con un estimador intervalar una probabilidad no nula de estar en lo correcto, en efecto,

$$\begin{aligned} P(\mu \in [\bar{X} - 1, \bar{X} + 1]) &= P(\bar{X} - 1 \leq \mu \leq \bar{X} + 1) = P(-1 \leq \mu - \bar{X} \leq 1) \\ &= P(-1 \leq \bar{X} - \mu \leq 1) = P\left(-\frac{1}{\sqrt{1/4}} \leq \frac{\bar{X} - \mu}{\sqrt{1/4}} \leq \frac{1}{\sqrt{1/4}}\right) \\ &= P\left(-2 \leq \frac{\bar{X} - \mu}{\sqrt{1/4}} \leq 2\right) = P(-2 \leq Z \leq 2) = 0.9544 \end{aligned}$$

en este caso $Z = (\bar{X} - \mu)/\sqrt{1/4} \sim N(0, 1)$.

INTERPRETACIÓN. De este modo, tenemos un 95 % de chances de cubrir el parámetro verdadero (desconocido) con nuestro estimador intervalar.

OBSERVACIÓN. En este contexto $P_\theta(\theta \in [L(\mathbf{x}), U(\mathbf{x})])$ se denomina *probabilidad de cobertura*

DEFINICIÓN 4.3. El *coeficiente de confianza* de $[L(\mathbf{x}), U(\mathbf{x})]$ es el ínfimo de las probabilidades de cobertura

$$\inf_{\theta} P_\theta(\theta \in [L(\mathbf{x}), U(\mathbf{x})])$$

OBSERVACIÓN. Estimadores intervalares en conjunto con una medida de confianza (coeficiente de confianza) son conocidos como *intervalos de confianza*.

A continuación consideraremos dos de los procedimientos más utilizados para construir intervalos de confianza.

4.1. Método de la Cantidad Pivotal

DEFINICIÓN 4.4. Una variable aleatoria $Q(\mathbf{X}; \theta) = Q(X_1, \dots, X_n; \theta)$ es una *cantidad pivotal* o pivote si la distribución de $Q(\mathbf{X}; \theta)$ **no** depende de θ . Esto es, si $\mathbf{X} \sim F(\mathbf{x}; \theta)$, entonces $Q(\mathbf{X}; \theta)$ tiene la misma distribución para todo valor de θ .

OBSERVACIÓN. La técnica confía en la habilidad de hallar un pivote y un conjunto A tal que el conjunto $\{\theta : Q(\mathbf{X}; \theta) \in A\}$ sea una estimación intervalar para θ .

EJEMPLO 4.5. Si X_1, \dots, X_n es una muestra aleatoria de tamaño n desde $N(\mu, \sigma^2)$, entonces

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1),$$

y por tanto es un pivote para μ (siempre que σ^2 sea conocido). Para cualquier constante a sigue que:

$$\begin{aligned} P\left(-a \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq a\right) &= P\left(-a \frac{\sigma}{\sqrt{n}} \leq \bar{X} - \mu \leq a \frac{\sigma}{\sqrt{n}}\right) \\ &= P\left(\bar{X} - a \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + a \frac{\sigma}{\sqrt{n}}\right), \end{aligned}$$

es decir obtenemos el intervalo de confianza

$$\left[\bar{X} - a \frac{\sigma}{\sqrt{n}}; \bar{X} + a \frac{\sigma}{\sqrt{n}}\right],$$

o bien

$$\left\{\mu : \bar{X} - a \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + a \frac{\sigma}{\sqrt{n}}\right\}.$$

Además suponga que $a = z_{1-\alpha/2}$ para un valor de α dado. Entonces, es fácil notar que

$$P\left(\mu \in \left[\bar{X} - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}; \bar{X} + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right]\right) = 1 - \alpha,$$

corresponde a un intervalo de confianza del $100(1 - \alpha)\%$ para μ .

Para el caso en que σ^2 sea desconocido podemos usar el pivote

$$T = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t(n-1),$$

es decir,

$$P(-a \leq T \leq a) = P\left(-a \leq \frac{\bar{X} - \mu}{s/\sqrt{n}} \leq a\right)$$

que lleva al intervalo de confianza

$$\left\{\mu : \bar{X} - t_{1-\alpha/2}(n-1) \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{1-\alpha/2}(n-1) \frac{s}{\sqrt{n}}\right\}.$$

OBSERVACIÓN. Note que los intervalos anteriores son simétricos.

Considere ahora,

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

que es cantidad pivotal y elija a y b , satisfaciendo que

$$\mathbf{P}(a \leq \chi^2 \leq b) = \mathbf{P}\left(a \leq \frac{(n-1)S^2}{\sigma^2} \leq b\right) = 1 - \alpha,$$

desde donde obtenemos

$$\left\{\sigma^2 : \frac{(n-1)S^2}{b} \leq \sigma^2 \leq \frac{(n-1)S^2}{a}\right\}.$$

Las elecciones de a y b que producen el intervalo con el coeficiente de confianza requerido son $a = \chi_{1-\alpha/2}^2(n-1)$ y $b = \chi_{\alpha/2}^2(n-1)$.

4.2. Intervalos de Confianza Asintóticos

DEFINICIÓN 4.6. Considere $\mathbf{SE} = \sqrt{\text{var}(\hat{\theta}_n)}$. Entonces $\widehat{\mathbf{SE}} = \sqrt{1/\mathcal{F}_n(\hat{\theta}_n)}$, luego un intervalo de confianza del $100(1 - \alpha)\%$ para θ es dado por¹

$$IC_n(\theta) = [\hat{\theta}_n - z_{1-\alpha/2}\widehat{\mathbf{SE}}, \hat{\theta}_n + z_{1-\alpha/2}\widehat{\mathbf{SE}}].$$

EJEMPLO 4.7. Sea X_1, \dots, X_n muestra aleatoria desde $\text{Ber}(p)$. Sabemos que el MLE de p es $\hat{p}_n = \frac{1}{n} \sum_{i=1}^n x_i$, y

$$\log f(x; p) = x \log p + (1 - x) \log(1 - p),$$

así

$$U(x; p) = \frac{x}{p} - \frac{1-x}{1-p}, \quad U'(x; p) = \frac{x}{p^2} + \frac{1-x}{(1-p)^2}.$$

De este modo,

$$\mathcal{F}_1(p) = \mathbf{E}\{-U'(X; p)\} = \frac{p}{p^2} + \frac{1-p}{(1-p)^2} = \frac{1}{p} + \frac{1}{1-p} = \frac{1}{p(1-p)},$$

de ahí que

$$\widehat{\mathbf{SE}} = \frac{1}{\sqrt{\mathcal{F}_n(\hat{p}_n)}} = \frac{1}{\sqrt{n\mathcal{F}_1(\hat{p}_n)}} = \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}},$$

luego, un intervalo de confianza del $100(1 - \alpha)\%$ para p es dado por

$$\hat{p}_n \mp z_{1-\alpha/2} \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}}.$$

OBSERVACIÓN. Considere $\lambda = g(\theta)$. Sabemos que el estimador ML de λ es dado por $\hat{\lambda}_n = g(\hat{\theta}_n)$. Además, usando el método Delta, sigue que

$$\frac{(\hat{\lambda}_n - \lambda)}{\widehat{\mathbf{SE}}(\hat{\lambda}_n)} \xrightarrow{D} \mathbf{N}(0, 1),$$

donde

$$\widehat{\mathbf{SE}}(\hat{\lambda}_n) = |g'(\hat{\theta}_n)| \widehat{\mathbf{SE}}(\hat{\theta}_n).$$

Lo que lleva al intervalo de confianza asintótico

$$IC_n(\lambda) = [\hat{\lambda}_n - z_{1-\alpha/2} \widehat{\mathbf{SE}}(\hat{\lambda}_n), \hat{\lambda}_n + z_{1-\alpha/2} \widehat{\mathbf{SE}}(\hat{\lambda}_n)].$$

¹En efecto, $\mathbf{P}_\theta(\theta \in IC_n(\theta)) \rightarrow 1 - \alpha$ para $n \rightarrow \infty$.

EJEMPLO 4.8. Considere X_1, \dots, X_n variables aleatorias IID desde una FE 1-paramétrica, y sea $\phi = \eta(\theta)$, $\gamma(\phi) = b(\theta)$. De este modo,

$$f(x; \phi) = \exp[\phi T(x) - \gamma(\phi)] h(x),$$

es decir, la log-verosimilitud para una única observación es dada por

$$\log f(x; \phi) = \phi T(x) - \gamma(\phi) + \log h(x).$$

Lo que lleva a,

$$U(x; \phi) = T(x) - \gamma'(\phi), \quad U'(x; \phi) = -\gamma''(\phi),$$

de ahí que, la información de Fisher es dada por

$$\mathcal{F}_1(\phi) = \mathbf{E}\{-U'(X; \phi)\} = \gamma''(\phi) = \text{var}(T(X)).$$

Es decir, el error estándar adopta la forma

$$\widehat{\text{SE}} = \frac{1}{\sqrt{\mathcal{F}_n(\hat{\phi}_n)}} = \frac{1}{\sqrt{n\mathcal{F}_1(\hat{\phi}_n)}} = \frac{1}{\sqrt{n\gamma''(\hat{\phi}_n)}},$$

donde $\hat{\phi}_n$ denota el MLE de ϕ (ver Ejemplo 3.16). Finalmente, un intervalo de confianza del $100(1 - \alpha)\%$ para ϕ es dado por

$$IC_n(\phi) = \left[\hat{\phi}_n \mp z_{1-\alpha/2} \frac{1}{\sqrt{n\gamma''(\hat{\phi}_n)}} \right].$$

Evidentemente, también podemos considerar un intervalo de confianza asintótico para $\theta = g(\phi) = \eta^{-1}(\phi)$ usando el método Delta.

4.3. Regiones de Confianza Asintóticas

Cuando θ es un vector k -dimensional, podemos definir una región de confianza (asintótica), mediante

$$\lim_{n \rightarrow \infty} P_\theta(\theta \in RC_n(\theta)) = 1 - \alpha,$$

cuando θ es el verdadero vector de parámetros. El mecanismo usado para construir una región de confianza asintótica está basada en el siguiente resultado.

RESULTADO 4.9. Sea $\{\mathbf{T}_n\}$ una secuencia de vectores aleatorios k -dimensionales tal que $\sqrt{n}(\mathbf{T}_n - \theta) \xrightarrow{D} N_k(\mathbf{0}, \Sigma)$ y sea $\{\mathbf{A}_n\}$ una secuencia de matrices aleatorias tal que $\mathbf{A}_n \xrightarrow{P} \mathbf{A}$, donde $\mathbf{A}\Sigma\mathbf{A} = \mathbf{A}$. Entonces

$$Q_n = n(\mathbf{T}_n - \theta)^\top \mathbf{A}_n(\mathbf{T}_n - \theta) \xrightarrow{D} \chi^2(k).$$

DEMOSTRACIÓN. Disponible en [Sen y Singer \(1993\)](#), página 137. □

Basado en el Resultado 3.40 y usando que $\mathcal{F}_1(\hat{\theta}_n) \xrightarrow{P} \mathcal{F}_1(\theta)$, tenemos

$$n(\hat{\theta}_n - \theta)^\top \mathcal{F}_1(\hat{\theta}_n)(\hat{\theta}_n - \theta) \xrightarrow{D} \chi^2(k).$$

Es decir, podemos construir una región de confianza del $100(1 - \alpha)\%$ para θ como:

$$RC_n(\theta) = \{\theta : n(\hat{\theta}_n - \theta)^\top \mathcal{F}_1(\hat{\theta}_n)(\hat{\theta}_n - \theta) \leq \chi_{1-\alpha}^2(k)\},$$

donde $\chi_{1-\alpha}^2(k)$ denota un valor cuantil $1 - \alpha$ de la distribución chi-cuadrado con k grados de libertad.

EJEMPLO 4.10. Considere X_1, \dots, X_n una muestra aleatoria desde $\mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Es fácil notar que $\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$ y

$$\mathbf{S}_n = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}_n)(\mathbf{X}_i - \bar{\mathbf{X}}_n)^\top \xrightarrow{P} \boldsymbol{\Sigma}.$$

De este modo, la estadística T^2 de Hotelling, satisface que

$$T_n^2 = n(\bar{\mathbf{X}}_n - \boldsymbol{\mu})^\top \mathbf{S}_n^{-1}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} \chi^2(p).$$

Luego, una región de confianza asintótica del $100(1 - \alpha)\%$ para $\boldsymbol{\mu}$ es dada por:

$$RC_n(\boldsymbol{\mu}) = \{\boldsymbol{\mu} : n(\bar{\mathbf{X}}_n - \boldsymbol{\mu})^\top \mathbf{S}_n^{-1}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \leq \chi_{1-\alpha}^2(p)\}.$$

Bibliografía

- Bates, D.M., Watts, D.G. (1981). A relative offset orthogonality convergence criterion for nonlinear least squares. *Technometrics* **23**, 179-183.
- Berndt, E.K., Hall, B.H., Hall, R.E., Hausman, J.A. (1974). Estimation and inference in nonlinear structural models. *Annals of Economic and Social Measurements* **3**, 653-665.
- Crudu, F., Osorio, F. (2020). Bilinear form test statistics for extremum estimation. *Economics Letters* **187**, 108885.
- Casella, G., Berger, R.L. (2002). *Statistical Inference (2nd Ed.)*. Duxbury, Pacific Grove.
- Gómez, E., Gómez-Villegas, M.A., Marín, J.M. (1988). A multivariate generalization of the power exponential family of distributions. *Communications in Statistics - Theory and Methods* **27**, 589-600.
- Jørgensen, B. (1997). *The Theory of Dispersion Models*. Chapman & Hall, London.
- Jørgensen, B., Labouriau, R. (1994). *Exponential Families and Theoretical Inference*. Lecture Notes, Department of Statistics, University of British Columbia.
- Lange, K., Sinsheimer, J.S. (1993). Normal/Independent distributions and their applications in robust regression. *Journal of Computational and Graphical Statistics* **2**, 175-198.
- Lange, K., Little, R.J.A., Taylor, J.M.G. (1989). Robust statistical modeling using the t distribution. *Journal of the American Statistical Association* **84**, 881-896.
- Lehmann, E.L. (1986). *Testing Statistical Hypotheses*. Wiley, New York.
- Mood, A.M., Graybill, F.A., Boes, D.C. (1974). *Introduction to the Theory of Statistics*, 3rd Edition. McGraw-Hill, New York.
- Osorio, F. (2019). heavy: Robust estimation using heavy-tailed distributions. *R package version 0.38.196*, URL: <https://CRAN.R-project.org/package=heavy>
- Pawitan, Y. (2001). *In All Likelihood: Statistical Modelling and Inference using Likelihood*. Oxford University Press.
- Pinheiro, J., Bates, D.M., DebRoy, S., Sarkar, D. (2019). nlme: Linear and nonlinear mixed effects models. *R package version 3.1-143*, URL: <https://CRAN.R-project.org/package=nlme>
- Rohde, C.A. (2014). *Introductory Statistical Inference with the Likelihood Function*. Springer, New York.
- Sen, P.K., Singer, J.M. (1993). *Large Sample Methods in Statistics*. Chapman & Hall, London.
- Wasserman, L. (2003). *All of Statistics: A concise course in statistical inference*. Springer, New York.