

NOTAS DE CLASE :

Análisis de Regresión

Felipe Osorio

DEPARTAMENTO DE MATEMÁTICA, UNIVERSIDAD TÉCNICA FEDERICO SANTA MARÍA

Índice general

Prefacio	v
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. Distribución normal multivariada	9
1.6. Alternativas a la distribución normal multivariada	16
1.7. Algunas distribuciones no centrales	21
1.8. Distribución de formas cuadráticas	24
Ejercicios	31
Capítulo 2. Inferencia en el Modelo Lineal	35
2.1. Definición de un modelo lineal	35
2.2. Estimación de parámetros en el modelo de regresión lineal	37
2.3. Aspectos numéricos de estimación LS en regresión lineal	41
2.4. Estimación bajo restricciones lineales	49
2.5. Test de hipótesis lineales	54
2.6. Regiones de confianza	58
Ejercicios	59
Capítulo 3. Chequeo del Modelo y Alternativas a Mínimos Cuadrados	61
3.1. Colinealidad	61
3.2. Errores correlacionados	69
3.3. Transformaciones estabilizadoras de varianza	73
3.4. Análisis de residuos y leverages	77
3.5. Diagnóstico de por eliminación de casos	79
3.6. Procedimientos para estimación robusta	83
Apéndice A. Elementos de Álgebra Matricial	97
A.1. Vectores y matrices	97
A.2. Definiciones básicas y propiedades	97
A.3. Inversa generalizada y sistemas de ecuaciones lineales	108
Apéndice B. Diferenciación matricial	111
B.1. Aproximación de primer orden	111
B.2. Funciones matriciales	112
B.3. Matriz Hessiana	114

B.4. Reglas fundamentales	115
Bibliografía	117

Prefacio

Estas notas de clase están asociadas a los contenidos de la asignatura *MAT-266: Análisis de Regresión*, dictada en el programa de Ingeniería Civil Matemática de la Universidad Técnica Federico Santa María. Aunque el documento se encuentra en una etapa bastante preliminar, espero ir puliendo el mismo para que se pueda convertir en un apunte que sirva de apoyo para los estudiantes de la asignatura.

Las notas se encuentran divididas en tres partes, con preliminares conteniendo resultados de la distribución normal y su conexión con formas cuadráticas. Luego, se presenta la inferencia en el modelo de regresión lineal y posteriormente hay una serie de resultados que permiten la crítica del proceso de modelación así como procedimientos alternativos en caso de que algunos de los supuestos básicos no sean satisfechos. El objetivo de este texto es proveer de una introducción rigurosa al tópico de regresión presentando también aplicaciones prácticas así como destacar los elementos necesarios para la implementación computacional de tales técnicas.

Agradezco al profesor Manuel Galea por haberme introducido en este tópico así como por su constante apoyo durante toda mi carrera académica. Adicionalmente, debo destacar el apoyo de los estudiantes que han participado de alguna versión de este curso, pues producto de sus comentarios y sugerencias este documento se ha visto notablemente mejorado.

Felipe Osorio
Valparaíso, Agosto 2021.

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 g de \mathbf{X} está dada por

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

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. Distribución normal multivariada

La distribución normal multivariada ocupa un rol central en inferencia multivariada así como en modelación estadística. En esta sección introducimos la distribución normal multivariada mediante tres definiciones equivalentes.

Una variable aleatoria (uni-dimensional) Z tiene una distribución normal con media cero y varianza uno si su función de densidad es de la forma

$$f(z) = (2\pi)^{-1/2} \exp\left(-\frac{1}{2}z^2\right), \quad z \in \mathbb{R},$$

en cuyo caso escribimos $Z \sim \mathbf{N}(0, 1)$. Más generalmente una variable aleatoria $Y \in \mathbb{R}$ tiene distribución normal con media $\mu \in \mathbb{R}$ y varianza $\sigma^2 \geq 0$ si

$$Y \stackrel{d}{=} \mu + \sigma Z, \quad Z \sim \mathbf{N}(0, 1),$$

en cuyo caso escribimos $Y \sim \mathbf{N}(\mu, \sigma^2)$. Cuando $\sigma^2 = 0$, la distribución $\mathbf{N}(\mu, \sigma^2)$ se interpreta como una distribución degenerada en μ . Si $Y \sim \mathbf{N}(\mu, \sigma^2)$, entonces su función característica adopta la forma

$$\varphi(t) = \exp\left(it\mu - \frac{1}{2}\sigma^2 t^2\right), \quad t \in \mathbb{R}.$$

Sea Z_1, \dots, Z_n variables aleatorias independientes cada una con distribución $N(0, 1)$ y considere el vector aleatorio $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$. De este modo, la densidad conjunta de \mathbf{Z} es dada por

$$\begin{aligned} f(\mathbf{z}) &= \prod_{i=1}^n (2\pi)^{-1/2} \exp\left(-\frac{1}{2}z_i^2\right) = (2\pi)^{-n/2} \exp\left(-\frac{1}{2}\sum_{i=1}^n z_i^2\right) \\ &= (2\pi)^{-n/2} \exp\left(-\frac{1}{2}\|\mathbf{z}\|^2\right), \end{aligned}$$

y anotamos $\mathbf{Z} \sim N_n(\mathbf{0}, \mathbf{I})$.

DEFINICIÓN 1.18. Un vector aleatorio p -dimensional, \mathbf{X} tiene distribución normal con vector de medias $\boldsymbol{\mu} \in \mathbb{R}^p$ y matriz de covarianza $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma} \geq \mathbf{0}$ sólo si, para todo vector \mathbf{t} la variable aleatoria (uni-dimensional) $\mathbf{t}^\top \mathbf{X}$ es normal univariada, en cuyo caso escribimos $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

OBSERVACIÓN 1.19. Note que en la definición anterior *no* se ha hecho supuestos respecto de la independencia de los componentes de \mathbf{X} .

RESULTADO 1.20. Suponga que $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 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 $\mathbf{Y} \sim N_m(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$.

DEMOSTRACIÓN. Sea $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$ y simplemente note que

$$\mathbf{t}^\top \mathbf{Y} = \mathbf{t}^\top \mathbf{A}\mathbf{X} + \mathbf{t}^\top \mathbf{b} = (\mathbf{A}^\top \mathbf{t})^\top \mathbf{X} + \mathbf{t}^\top \mathbf{b} = \mathbf{h}^\top \mathbf{X} + c,$$

por la Definición 1.18 tenemos que $\mathbf{h}^\top \mathbf{X}$ es normal y como c es una constante, sigue que $\mathbf{t}^\top \mathbf{Y}$ tiene distribución normal multivariada. \square

A partir del resultado anterior sigue que todas las distribuciones marginales de \mathbf{X} también son normalmente distribuidas. En particular, también permite apreciar que la distribución normal satisface la siguiente propiedad relativa a la simetría multivariada:

$$\mathbf{Z} \sim N_p(\mathbf{0}, \sigma^2 \mathbf{I}_p) \implies \mathbf{Q}\mathbf{Z} \stackrel{d}{=} \mathbf{Z}, \quad \forall \mathbf{Q} \in \mathcal{O}_p.$$

RESULTADO 1.21. Si $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, entonces la función característica de \mathbf{X} es dada por

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

DEMOSTRACIÓN. Sabemos que la función característica de un vector aleatorio, satisface

$$\varphi_X(\mathbf{t}) = \mathbb{E}\{\exp(i\mathbf{t}^\top \mathbf{X})\} = \varphi_{\mathbf{t}^\top \mathbf{X}}(1),$$

donde la función característica de la variable aleatoria uni-dimensional $Y = \mathbf{t}^\top \mathbf{X}$ es evaluada en 1. Como $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ sólo si $\mathbf{t}^\top \mathbf{X} \sim N_1(\mathbf{t}^\top \boldsymbol{\mu}, \mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t})$, tenemos

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

En efecto, sea $\boldsymbol{\Sigma}$ matriz de covarianza $p \times p$ semidefinida positiva de rango r y sea Z_1, \dots, Z_r variables aleatorias IID $N(0, 1)$. Entonces el vector $\mathbf{Z} = (Z_1, \dots, Z_r)^\top$ tiene función característica

$$\begin{aligned} \varphi_Z(\mathbf{t}) &= \mathbb{E}\{\exp(i\mathbf{t}^\top \mathbf{Z})\} = \prod_{j=1}^r \mathbb{E}\{\exp(it_j Z_j)\} \\ &= \prod_{j=1}^r \exp\left(-\frac{1}{2}t_j^2\right) = \exp\left(-\frac{1}{2}\mathbf{t}^\top \mathbf{t}\right). \end{aligned}$$

Considere

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{B}\mathbf{Z},$$

donde $\mathbf{B} \in \mathbb{R}^{p \times r}$ con $\text{rg}(\mathbf{B}) = r$, tal que $\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}^\top$ y $\boldsymbol{\mu} \in \mathbb{R}^p$. De este modo, \mathbf{X} tiene función característica

$$\begin{aligned} \varphi_{\mathbf{X}}(\mathbf{t}) &= \mathbb{E}\{\exp(i\mathbf{t}^\top \mathbf{X})\} = \mathbb{E}\{\exp(i\mathbf{t}^\top (\boldsymbol{\mu} + \mathbf{B}\mathbf{Z}))\} \\ &= \exp(i\mathbf{t}^\top \boldsymbol{\mu}) \mathbb{E}\{\exp(i\mathbf{t}^\top \mathbf{B}\mathbf{Z})\} = \exp(i\mathbf{t}^\top \boldsymbol{\mu}) \varphi_{\mathbf{Z}}(\mathbf{h}), \quad \mathbf{h} = \mathbf{B}^\top \mathbf{t} \\ &= \exp(i\mathbf{t}^\top \boldsymbol{\mu}) \exp(-\frac{1}{2}\mathbf{t}^\top \mathbf{B}\mathbf{B}^\top \mathbf{t}) = \exp(i\mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}). \end{aligned}$$

□

OBSERVACIÓN 1.22. El Resultado 1.20 puede ser demostrado de manera bastante simple usando la función característica (ver Ejercicio 1.3).

RESULTADO 1.23. Si $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$. Entonces

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

DEMOSTRACIÓN. Para mostrar el resultado deseado, podemos calcular el primer y segundo diferencial de la función característica del vector aleatorio $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$. Debemos calcular,

$$d\varphi_{\mathbf{Z}}(\mathbf{t}) = -\varphi_{\mathbf{Z}}(\mathbf{t})\mathbf{t}^\top d\mathbf{t},$$

y

$$\begin{aligned} d^2 \varphi_{\mathbf{Z}}(\mathbf{t}) &= -d\varphi_{\mathbf{Z}}(\mathbf{t})\mathbf{t}^\top d\mathbf{t} - \varphi_{\mathbf{Z}}(\mathbf{t})(d\mathbf{t})^\top d\mathbf{t} \\ &= \varphi_{\mathbf{Z}}(\mathbf{t})(d\mathbf{t})^\top \mathbf{t}\mathbf{t}^\top d\mathbf{t} - \varphi_{\mathbf{Z}}(\mathbf{t})(d\mathbf{t})^\top d\mathbf{t} \\ &= \varphi_{\mathbf{Z}}(\mathbf{t})(d\mathbf{t})^\top (\mathbf{t}\mathbf{t}^\top - \mathbf{I}) d\mathbf{t}, \end{aligned}$$

de ahí que

$$\frac{\partial \varphi_{\mathbf{Z}}(\mathbf{t})}{\partial \mathbf{t}} = -\varphi_{\mathbf{Z}}(\mathbf{t})\mathbf{t}, \quad \frac{\partial^2 \varphi_{\mathbf{Z}}(\mathbf{t})}{\partial \mathbf{t} \partial \mathbf{t}^\top} = \varphi_{\mathbf{Z}}(\mathbf{t})(\mathbf{t}\mathbf{t}^\top - \mathbf{I}).$$

Ahora, el vector de medias y matriz de covarianzas están dadas por

$$\begin{aligned} \mathbb{E}(\mathbf{Z}) &= i^{-1} \frac{\partial \varphi_{\mathbf{Z}}(\mathbf{t})}{\partial \mathbf{t}} \Big|_{\mathbf{t}=\mathbf{0}} = \mathbf{0}, \\ \mathbb{E}(\mathbf{Z}\mathbf{Z}^\top) &= i^{-2} \frac{\partial^2 \varphi_{\mathbf{Z}}(\mathbf{t})}{\partial \mathbf{t} \partial \mathbf{t}^\top} \Big|_{\mathbf{t}=\mathbf{0}} = \mathbf{I} = \text{Cov}(\mathbf{Z}). \end{aligned}$$

□

OBSERVACIÓN 1.24. Considere

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{B}\mathbf{Z}, \quad \boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}^\top,$$

con $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$. Usando los Resultados 1.20 y 1.23, sigue que

$$\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu} + \mathbf{B}\mathbb{E}(\mathbf{Z}) = \boldsymbol{\mu}, \quad \text{Cov}(\mathbf{X}) = \mathbf{B}\text{Cov}(\mathbf{Z})\mathbf{B}^\top = \boldsymbol{\Sigma}.$$

RESULTADO 1.25. Si $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, entonces la distribución marginal de cualquier subconjunto de k ($< p$) componentes de \mathbf{X} es normal k -variada.

DEMOSTRACIÓN. Considere la siguiente partición:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}, \quad (1.3)$$

donde \mathbf{X}_1 y $\boldsymbol{\mu}_1$ son vectores $k \times 1$ y $\boldsymbol{\Sigma}_{11}$ es $k \times k$. Aplicando el Resultado 1.20 con

$$\mathbf{A} = (\mathbf{I}_k, \mathbf{0}) \in \mathbb{R}^{k \times p} \quad \text{y} \quad \mathbf{b} = \mathbf{0},$$

sigue inmediatamente que $\mathbf{X}_1 \sim \mathcal{N}_k(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$. \square

Una consecuencia de este resultado es que la distribución marginal de *cada* componente de \mathbf{X} es normal univariada.

OBSERVACIÓN 1.26. La inversa del Resultado 1.25 *no* es verdad en general. Es decir, que cada componente de un vector aleatorio tenga distribución normal no implica que todo el vector siga una distribución normal multivariada.

RESULTADO 1.27. Si $\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ y \mathbf{X} , $\boldsymbol{\mu}$ y $\boldsymbol{\Sigma}$ son particionadas como en la Ecuación (1.3). Entonces los vectores \mathbf{X}_1 y \mathbf{X}_2 son independientes si y sólo si $\boldsymbol{\Sigma}_{12} = \mathbf{0}$.

DEMOSTRACIÓN. Note que $\text{Cov}(\mathbf{X}_1, \mathbf{X}_2) = \boldsymbol{\Sigma}_{12}$, así la independencia entre \mathbf{X}_1 y \mathbf{X}_2 implica que $\boldsymbol{\Sigma}_{12} = \mathbf{0}$. Suponga ahora que $\boldsymbol{\Sigma}_{12} = \mathbf{0}$. Entonces la función característica

$$\begin{aligned} \varphi_{\mathbf{X}}(\mathbf{t}) &= \exp(i\mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}) \\ &= \exp(it_1^\top \boldsymbol{\mu}_1 + it_2^\top \boldsymbol{\mu}_2 - \frac{1}{2}t_1^\top \boldsymbol{\Sigma}_{11} t_1 - \frac{1}{2}t_2^\top \boldsymbol{\Sigma}_{22} t_2) \\ &= \exp(it_1^\top \boldsymbol{\mu}_1 - \frac{1}{2}t_1^\top \boldsymbol{\Sigma}_{11} t_1) \exp(it_2^\top \boldsymbol{\mu}_2 - \frac{1}{2}t_2^\top \boldsymbol{\Sigma}_{22} t_2) \\ &= \varphi_{\mathbf{X}_1}(\mathbf{t}_1) \varphi_{\mathbf{X}_2}(\mathbf{t}_2), \end{aligned}$$

es decir, $\mathbf{X}_1 \sim \mathcal{N}_k(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$ es independiente de $\mathbf{X}_2 \sim \mathcal{N}_{p-k}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$. \square

DEFINICIÓN 1.28. Si $\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ y $\boldsymbol{\Sigma}$ es definida positiva, entonces la densidad de \mathbf{X} asume la forma

$$f_{\mathbf{X}}(\mathbf{x}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\}, \quad \mathbf{x} \in \mathbb{R}^p.$$

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

$$f_{\mathbf{Z}}(\mathbf{z}) = (2\pi)^{-p/2} \exp(-\frac{1}{2}\|\mathbf{z}\|^2).$$

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

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

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

EJEMPLO 1.29. Sea $\mathbf{X} \sim N_2(\mathbf{0}, \Sigma)$ donde

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

En cuyo caso, la función de densidad es dada por:

$$f(\mathbf{x}) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)}(x_1^2 + x_2^2 - 2\rho x_1 x_2) \right\}.$$

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

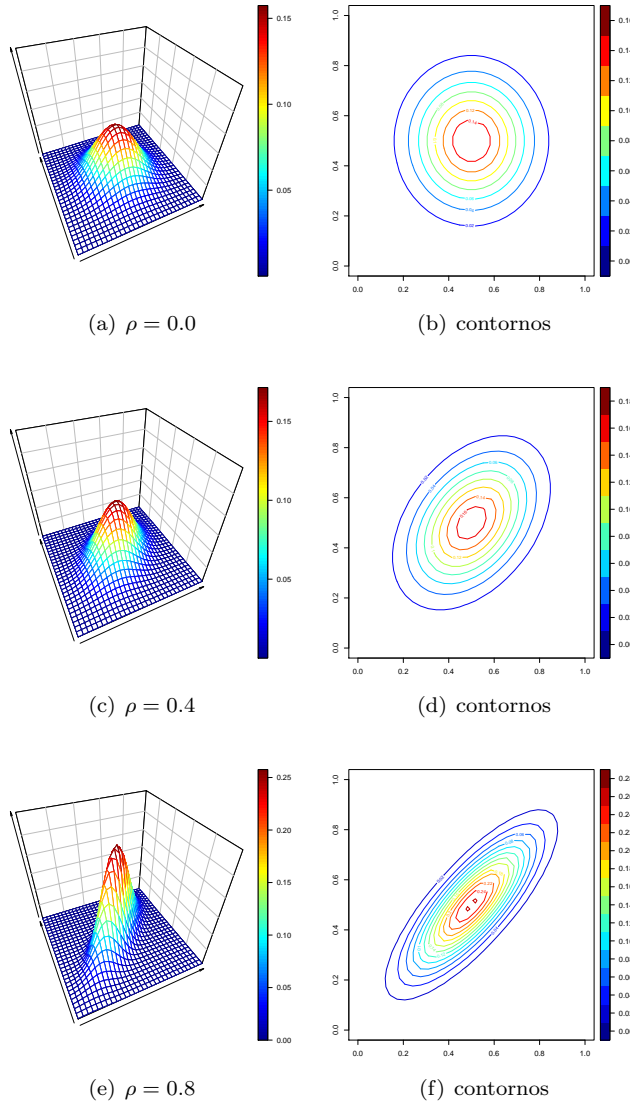


FIGURA 1. Densidad de $\mathbf{X} \sim N_2(\mathbf{0}, \Sigma)$ para $\rho = 0.0, 0.4$ y 0.8 .

Es fácil apreciar que la función de densidad es constante sobre el elipsoide

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

en \mathbb{R}^p para todo $\lambda > 0$. Este elipsoide tiene centro $\boldsymbol{\mu}$, mientras que $\boldsymbol{\Sigma}$ determina su forma y orientación. Además, la variable aleatoria

$$(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) = \mathbf{Z}^\top \mathbf{Z} = \sum_{i=1}^p Z_i^2, \quad (1.4)$$

sigue una distribución chi-cuadrado con p grados de libertad y la cantidad $D = \{(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu})\}^{1/2}$ se conoce como *distancia de Mahalanobis* de \mathbf{X} a $\boldsymbol{\mu}$.

OBSERVACIÓN 1.30. Para la existencia de densidad hemos asumido que $\boldsymbol{\Sigma} > \mathbf{0}$. En el caso de que $\boldsymbol{\Sigma} \geq \mathbf{0}$ decimos que \mathbf{X} sigue una distribución normal singular.

Para introducir una definición de la función de densidad asociada a una variable con distribución normal singular, note que $X \sim \mathcal{N}(\mu, \sigma^2)$ con $\sigma^2 = 0 \Leftrightarrow x = \mu$ con probabilidad 1 (pues si $\sigma^2 = 0$, $P(X = \mu) = \lim_{n \rightarrow \infty} P(|X - \mu| < 1/n) = 0$, $\forall n$).

Considere $\mathbf{Y} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ con $\text{rg}(\boldsymbol{\Sigma}) = r < p$. Entonces, podemos escribir

$$\boldsymbol{\Sigma} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top = (\mathbf{U}_1, \mathbf{U}_2) \begin{pmatrix} \boldsymbol{\Lambda}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{U}_1^\top \\ \mathbf{U}_2^\top \end{pmatrix} = \mathbf{U}_1 \boldsymbol{\Lambda}_1 \mathbf{U}_1^\top,$$

donde $\boldsymbol{\Lambda}_1 = \text{diag}(\lambda_1, \dots, \lambda_r)$. De este modo, es claro que

$$\mathbf{U}^\top \boldsymbol{\Sigma} \mathbf{U} \implies \mathbf{U}_2^\top \boldsymbol{\Sigma} \mathbf{U}_2 = \mathbf{0},$$

es decir, tenemos que $\mathbf{U}_2^\top (\mathbf{Y} - \boldsymbol{\mu}) = \mathbf{0}$ con probabilidad 1. Mientras que

$$\mathbf{U}_1^\top (\mathbf{Y} - \boldsymbol{\mu}) \sim \mathcal{N}_r(\mathbf{0}, \boldsymbol{\Lambda}_1).$$

Además $\boldsymbol{\Sigma}^- = \mathbf{U}_1 \boldsymbol{\Lambda}_1^{-1} \mathbf{U}_1^\top = \mathbf{U}_1 (\mathbf{U}_1^\top \boldsymbol{\Sigma} \mathbf{U}_1)^{-1} \mathbf{U}_1^\top$. Así, \mathbf{Y} tiene la siguiente densidad normal (singular)

$$\begin{aligned} f_Y(\mathbf{y}) &= |2\pi \mathbf{U}_1^\top \boldsymbol{\Sigma} \mathbf{U}_1|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{U}_1^\top (\mathbf{y} - \boldsymbol{\mu}))^\top (\mathbf{U}_1^\top \boldsymbol{\Sigma} \mathbf{U}_1)^{-1} \mathbf{U}_1^\top (\mathbf{y} - \boldsymbol{\mu})\right\} \\ &= (2\pi)^{-r/2} |\boldsymbol{\Lambda}_1|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{U}_1 \boldsymbol{\Lambda}_1^{-1} \mathbf{U}_1^\top (\mathbf{y} - \boldsymbol{\mu})\right\}. \end{aligned}$$

El siguiente resultado presenta la distribución condicional de un vector aleatorio con distribución normal multivariada.

RESULTADO 1.31. Sea $\mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ y particione \mathbf{X} , $\boldsymbol{\mu}$ y $\boldsymbol{\Sigma}$ como:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix},$$

donde \mathbf{X}_1 y $\boldsymbol{\mu}_1$ son vectores $k \times 1$, mientras que $\boldsymbol{\Sigma}_{11}$ es matriz $k \times k$. Sea $\boldsymbol{\Sigma}_{22}^-$ una inversa generalizada de $\boldsymbol{\Sigma}_{22}$, esto es, una matriz que satisface

$$\boldsymbol{\Sigma}_{22} \boldsymbol{\Sigma}_{22}^- \boldsymbol{\Sigma}_{22} = \boldsymbol{\Sigma}_{22},$$

y sea $\boldsymbol{\Sigma}_{11.2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^- \boldsymbol{\Sigma}_{21}$. Entonces

- (a) $\mathbf{X}_1 - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^- \mathbf{X}_2 \sim \mathcal{N}_k(\boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^- \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{11.2})$ y es independiente de \mathbf{X}_2 .
- (b) La distribución condicional

$$(\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2) \sim \mathcal{N}_k(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^- (\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11.2}).$$

DEMOSTRACIÓN. Considere la transformación lineal

$$\mathbf{Y} = \begin{pmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{I}_k & -\mathbf{B} \\ \mathbf{0} & \mathbf{I}_{p-k} \end{pmatrix} \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} = \mathbf{C}\mathbf{X},$$

sigue que $\mathbf{Y} \sim \mathcal{N}_p(\mathbf{C}\boldsymbol{\mu}, \mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^\top)$, donde

$$\begin{aligned} \mathbf{C}\boldsymbol{\mu} &= \begin{pmatrix} \mathbf{I}_k & -\mathbf{B} \\ \mathbf{0} & \mathbf{I}_{p-k} \end{pmatrix} \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_1 - \mathbf{B}\boldsymbol{\mu}_2 \\ \boldsymbol{\mu}_2 \end{pmatrix} \\ \mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^\top &= \begin{pmatrix} \mathbf{I}_k & -\mathbf{B} \\ \mathbf{0} & \mathbf{I}_{p-k} \end{pmatrix} \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \\ -\mathbf{B}^\top & \mathbf{I}_{p-k} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\Sigma}_{11} - \mathbf{B}\boldsymbol{\Sigma}_{21} - \boldsymbol{\Sigma}_{12}\mathbf{B}^\top + \mathbf{B}\boldsymbol{\Sigma}_{22}\mathbf{B}^\top & \boldsymbol{\Sigma}_{12} - \mathbf{B}\boldsymbol{\Sigma}_{22} \\ \boldsymbol{\Sigma}_{21} - \boldsymbol{\Sigma}_{22}\mathbf{B}^\top & \boldsymbol{\Sigma}_{22} \end{pmatrix}. \end{aligned}$$

De este modo, nuestro interés es escoger $\boldsymbol{\Sigma}_{12} - \mathbf{B}\boldsymbol{\Sigma}_{22} = \mathbf{0}$. Es decir, $\boldsymbol{\Sigma}_{12} = \mathbf{B}\boldsymbol{\Sigma}_{22}$. Por otro lado, notando que

$$\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-\boldsymbol{\Sigma}_{22} = \mathbf{B}\boldsymbol{\Sigma}_{22}\boldsymbol{\Sigma}_{22}^-\boldsymbol{\Sigma}_{22} = \mathbf{B}\boldsymbol{\Sigma}_{22} = \boldsymbol{\Sigma}_{12},$$

sigue que $\boldsymbol{\Sigma}_{12}\mathbf{B}^\top = \mathbf{B}\boldsymbol{\Sigma}_{22}\mathbf{B}^\top$ (y análogamente $\mathbf{B}\boldsymbol{\Sigma}_{21} = \mathbf{B}\boldsymbol{\Sigma}_{22}\mathbf{B}^\top$). Esto es, si \mathbf{B} es escogida como $\mathbf{B} = \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-$, entonces \mathbf{Y}_1 y \mathbf{Y}_2 son independientes con distribución conjunta

$$\begin{pmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{X}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-\mathbf{X}_2 \\ \mathbf{X}_2 \end{pmatrix} \sim \mathcal{N}_p\left(\begin{pmatrix} \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-\boldsymbol{\mu}_2 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{11.2} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{22} \end{pmatrix}\right).$$

Esto muestra la parte (a). Para notar la parte (b), note que las densidades de \mathbf{Y}_1 y \mathbf{Y}_2 están dadas por

$$\begin{aligned} g(\mathbf{y}_1; \boldsymbol{\delta}_{1.2}, \boldsymbol{\Sigma}_{11.2}) &= |2\pi\boldsymbol{\Sigma}_{11.2}|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{y}_1 - \boldsymbol{\delta}_{1.2})^\top \boldsymbol{\Sigma}_{11.2}^{-1}(\mathbf{y}_1 - \boldsymbol{\delta}_{1.2})\right\} \\ f_2(\mathbf{y}_2; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}) &= |2\pi\boldsymbol{\Sigma}_{22}|^{-1/2} \exp\left\{-\frac{1}{2}(\mathbf{y}_2 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_{22}^{-1}(\mathbf{y}_2 - \boldsymbol{\mu}_2)\right\}, \end{aligned}$$

y la densidad conjunta para $\mathbf{Y} = (\mathbf{Y}_1^\top, \mathbf{Y}_2^\top)^\top$ adopta la forma

$$f(\mathbf{y}_1, \mathbf{y}_2; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = g(\mathbf{y}_1; \boldsymbol{\delta}_{1.2}, \boldsymbol{\Sigma}_{11.2}) f_2(\mathbf{y}_2; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}).$$

Como

$$f(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = f_{1|2}(\mathbf{x}_1; \boldsymbol{\mu}, \boldsymbol{\Sigma}|\mathbf{x}_2) f_2(\mathbf{x}_2; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}),$$

entonces, la densidad condicional de \mathbf{X}_1 dado $\mathbf{X}_2 = \mathbf{x}_2$ debe ser $g(\mathbf{y}_1; \boldsymbol{\delta}_{1.2}, \boldsymbol{\Sigma}_{11.2})$. Además, es fácil notar que la forma cuadrática

$$\begin{aligned} q(\mathbf{y}_1; \boldsymbol{\mu}_{1.2}, \boldsymbol{\Sigma}_{11.2}) &= (\mathbf{y}_1 - \boldsymbol{\delta}_{1.2})^\top \boldsymbol{\Sigma}_{11.2}^{-1}(\mathbf{y}_1 - \boldsymbol{\delta}_{1.2}) \\ &= (\mathbf{x}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-\mathbf{x}_2 - \boldsymbol{\delta}_{1.2})^\top \boldsymbol{\Sigma}_{11.2}^{-1}(\mathbf{x}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-\mathbf{x}_2 - \boldsymbol{\delta}_{1.2}) \\ &= (\mathbf{x}_1 - \boldsymbol{\mu}_{1.2})^\top \boldsymbol{\Sigma}_{11.2}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_{1.2}), \end{aligned}$$

donde

$$\boldsymbol{\mu}_{1.2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-(\mathbf{x}_2 - \boldsymbol{\mu}_2),$$

lo que muestra el resultado. \square

OBSERVACIÓN 1.32. La esperanza de la distribución condicional de \mathbf{X}_1 dado \mathbf{X}_2 , es decir

$$\mathbb{E}(\mathbf{X}_1|\mathbf{X}_2 = \mathbf{x}_2) = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-(\mathbf{x}_2 - \boldsymbol{\mu}_2),$$

se denomina *función de regresión* de \mathbf{X}_1 sobre \mathbf{X}_2 con coeficientes de regresión $\mathbf{B} = \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^-$. Esta es una función lineal de \mathbf{X}_2 y la matriz de covarianza $\boldsymbol{\Sigma}_{11.2}$ no depende de \mathbf{X}_2 .

RESULTADO 1.33. Sea $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ y considere $\mathbf{Y}_1 = \mathbf{A}_1 \mathbf{X}$, $\mathbf{Y}_2 = \mathbf{A}_2 \mathbf{X}$ dos funciones lineales del vector aleatorio \mathbf{X} . La covarianza entre \mathbf{Y}_1 y \mathbf{Y}_2 es dada por

$$\text{Cov}(\mathbf{Y}_1, \mathbf{Y}_2) = \mathbf{A}_1 \text{Cov}(\mathbf{X}, \mathbf{X}) \mathbf{A}_2^\top = \mathbf{A}_1 \boldsymbol{\Sigma} \mathbf{A}_2^\top$$

Este resultado permite obtener una condición para la independencia entre dos formas lineales en variables aleatorias normales, estos es \mathbf{Y}_1 y \mathbf{Y}_2 serán independientes si y sólo si $\mathbf{A}_1 \boldsymbol{\Sigma} \mathbf{A}_2^\top = \mathbf{0}$.

EJEMPLO 1.34. Considere X_1, \dots, X_n una muestra aleatoria desde $\mathbf{N}(\mu, \sigma^2)$ y sea $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$ el vector de datos centrados con $Z_i = X_i - \bar{X}$, $i = 1, \dots, n$, donde $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Podemos escribir

$$\bar{X} = \frac{1}{n} \mathbf{1}^\top \mathbf{X}, \quad \mathbf{Z} = \mathbf{C} \mathbf{X},$$

con $\mathbf{C} = \mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}^\top$ la matriz de centrado. Tenemos que $\mathbf{X} \sim \mathbf{N}_n(\mu \mathbf{1}, \sigma^2 \mathbf{I}_n)$ y \bar{X} con \mathbf{Z} son independientes pues $\mathbf{C} \mathbf{1} = \mathbf{0}$.

EJEMPLO 1.35. Sea $\mathbf{X} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$ y considere las transformaciones $\mathbf{Y}_1 = \mathbf{A} \mathbf{X}$ y $\mathbf{Y}_2 = (\mathbf{I} - \mathbf{A}^+ \mathbf{A})^\top \mathbf{X}$. De este modo

$$\text{Cov}(\mathbf{Y}_1, \mathbf{Y}_2) = \text{Cov}(\mathbf{A} \mathbf{X}, (\mathbf{I} - \mathbf{A}^+ \mathbf{A})^\top \mathbf{X}) = \sigma^2 \mathbf{A} (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) = \mathbf{0},$$

pues $\mathbf{A} \mathbf{A}^+ \mathbf{A} = \mathbf{A}$ y \mathbf{Y}_1 con \mathbf{Y}_2 son independientes.

1.6. Alternativas a la distribución normal multivariada

La distribución normal multivariada es de importancia fundamental en la teoría clásica de modelos lineales así como para análisis multivariado. A pesar de su uso amplio, es bien sabido que la inferencia estadística basada en la distribución normal es vulnerable a la presencia de datos atípicos, esto ha motivado considerar distribuciones alternativas que eviten este tipo de limitaciones. En esta dirección, varios autores han sugerido utilizar la clase de distribuciones elípticas (ver, por ejemplo, Fang et al., 1990; Arellano, 1994) particularmente debido al hecho de incluir distribuciones con colas más pesadas que la normal, tales como la t de Student, exponencial potencia y normal contaminada, entre otras. Una subclase importante de la familia de distribuciones elípticas es la clase de distribuciones de mezcla de escala normal (Andrews y Mallows, 1974) la que tiene propiedades similares a la distribución normal, es relativamente simple de trabajar y permite proponer procedimientos para estimación robusta. A continuación se presenta la definición y algunos ejemplos de distribuciones en la clase elíptica.

DEFINICIÓN 1.36. Sea \mathbf{U} vector aleatorio $p \times 1$ con *distribución uniforme* sobre el conjunto

$$\mathcal{S}_p = \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\| = 1\}, \quad (1.5)$$

esto es \mathcal{S}_p denota la *superficie de la esfera unitaria* en \mathbb{R}^p . En cuyo caso anotamos $\mathbf{U} \sim \mathbf{U}(\mathcal{S}_p)$.

PROPIEDAD 1.37. Si $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$, entonces $\mathbf{U} = (\mathbf{Z} / \|\mathbf{Z}\|)^\top \sim \mathbf{U}(\mathcal{S}_p)$, donde

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

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,

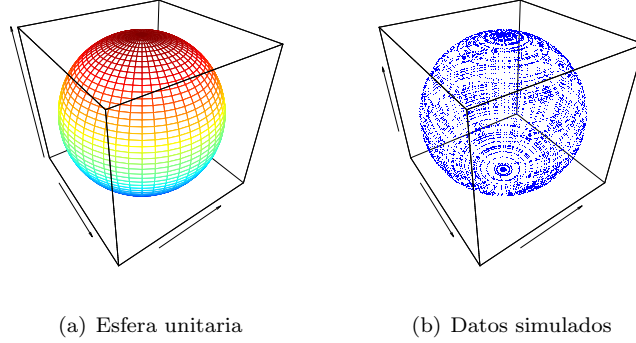


FIGURA 2. Esfera unitaria y datos generados sobre la superficie \mathcal{S}_p .

DEFINICIÓN 1.38. 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 1.39. Sea $\mathbf{U} \sim \mathcal{U}(\mathcal{S}_p)$, entonces es bastante obvio que $\mathbf{Q}\mathbf{U} \stackrel{d}{=} \mathbf{U}$.

EJEMPLO 1.40. Suponga $\mathbf{X} \sim \mathcal{N}_p(\mathbf{0}, \sigma^2 \mathbf{I})$. Tenemos que

$$\mathbf{Q}\mathbf{X} \sim \mathcal{N}_p(\mathbf{0}, \sigma^2 \mathbf{I}),$$

para $\mathbf{Q} \in \mathcal{O}_p$, es decir $\mathbf{Q}\mathbf{X} \stackrel{d}{=} \mathbf{X}$ tiene simetría esférica.

DEFINICIÓN 1.41. Un vector aleatorio p -dimensional tiene *distribución esférica* sólo si su función característica satisface

- (a) $\varphi(\mathbf{Q}^\top \mathbf{t}) = \varphi(\mathbf{t})$, para todo $\mathbf{Q} \in \mathcal{O}_p$.
- (b) Existe una función $\psi(\cdot)$ de una variable escalar tal que $\varphi(\mathbf{t}) = \psi(\mathbf{t}^\top \mathbf{t})$.

En este caso escribimos $\mathbf{X} \sim \mathcal{S}_p(\psi)$.

EJEMPLO 1.42. 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 1.43. Sea $\psi(\mathbf{t}^\top \mathbf{t})$ la función característica del vector aleatorio \mathbf{X} . Entonces \mathbf{X} tiene representación estocástica

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

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

RESULTADO 1.44. Suponga que $\mathbf{X} \stackrel{d}{=} \mathbf{R}\mathbf{U} \sim \mathcal{S}_p(\psi)$ ($P(\mathbf{X} = \mathbf{0}) = 0$), entonces

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

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

RESULTADO 1.45. *El vector de medias y la matriz de covarianza de $\mathbf{U} \sim \mathbf{U}(\mathcal{S}_p)$ son:*

$$\mathbf{E}(\mathbf{U}) = \mathbf{0}, \quad \text{Cov}(\mathbf{U}) = \frac{1}{p} \mathbf{I}_p,$$

respectivamente.

DEMOSTRACIÓN. Sea $\mathbf{X} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$, tenemos que $\mathbf{X} \stackrel{d}{=} \|\mathbf{X}\| \mathbf{U}$, con $\|\mathbf{X}\|$ independiente de \mathbf{U} . Sabemos que $\|\mathbf{X}\|^2 \sim \chi^2(p)$. Dado que

$$\mathbf{E}(\mathbf{X}) = \mathbf{0}, \quad \mathbf{E}(\|\mathbf{X}\|) > 0, \quad \text{y} \quad \mathbf{E}(\|\mathbf{X}\|^2) = p, \quad \text{Cov}(\mathbf{X}) = \mathbf{I}_p,$$

el resultado sigue. \square

RESULTADO 1.46. *Si $\mathbf{X} \stackrel{d}{=} R\mathbf{U} \sim \mathbf{S}_p(g)$ y $\mathbf{E}(R^2) < \infty$. Entonces,*

$$\mathbf{E}(\mathbf{X}) = \mathbf{0}, \quad \text{Cov}(\mathbf{X}) = \frac{\mathbf{E}(R^2)}{p} \mathbf{I}_p,$$

respectivamente.

DEMOSTRACIÓN. En efecto, como R y \mathbf{U} son independientes, sigue que

$$\mathbf{E}(\mathbf{X}) = \mathbf{E}(R) \mathbf{E}(\mathbf{U}) = \mathbf{0},$$

$$\text{Cov}(\mathbf{X}) = \mathbf{E}(R^2) \mathbf{E}(\mathbf{U}\mathbf{U}^\top) = \mathbf{E}(R^2) \text{Cov}(\mathbf{U}) = \frac{\mathbf{E}(R^2)}{p} \mathbf{I}_p,$$

siempre que $\mathbf{E}(R) < \infty$ y $\mathbf{E}(R^2) < \infty$. \square

DEFINICIÓN 1.47. Un vector aleatorio $p \times 1$, \mathbf{X} tiene *distribución de contornos elípticos* con parámetros $\boldsymbol{\mu} \in \mathbb{R}^p$ y $\boldsymbol{\Sigma} \geq \mathbf{0}$ si

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

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

OBSERVACIÓN 1.48. La función característica de $\mathbf{X} \sim \text{EC}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \psi)$ es de la forma

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

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 $\mathbf{B}\mathbf{B}^\top = \boldsymbol{\Sigma}$.

RESULTADO 1.49. *Suponga que $\mathbf{X} \sim \text{EC}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \psi)$ y $\mathbf{E}(R^2) < \infty$. Entonces*

$$\mathbf{E}(\mathbf{X}) = \boldsymbol{\mu}, \quad \text{Cov}(\mathbf{X}) = \frac{\mathbf{E}(R^2)}{p} \boldsymbol{\Sigma}.$$

DEFINICIÓN 1.50. 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:

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

y escribimos $\mathbf{X} \sim \text{EC}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; g)$.

OBSERVACIÓN 1.51. Asuma que $\mathbf{X} \sim \text{EC}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \psi)$ con $\text{rg}(\boldsymbol{\Sigma}) = k$. Entonces,

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

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

EJEMPLO 1.52. 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): $g(u) = c_4 \exp(-u^\lambda/2)$, $\lambda = 2$.

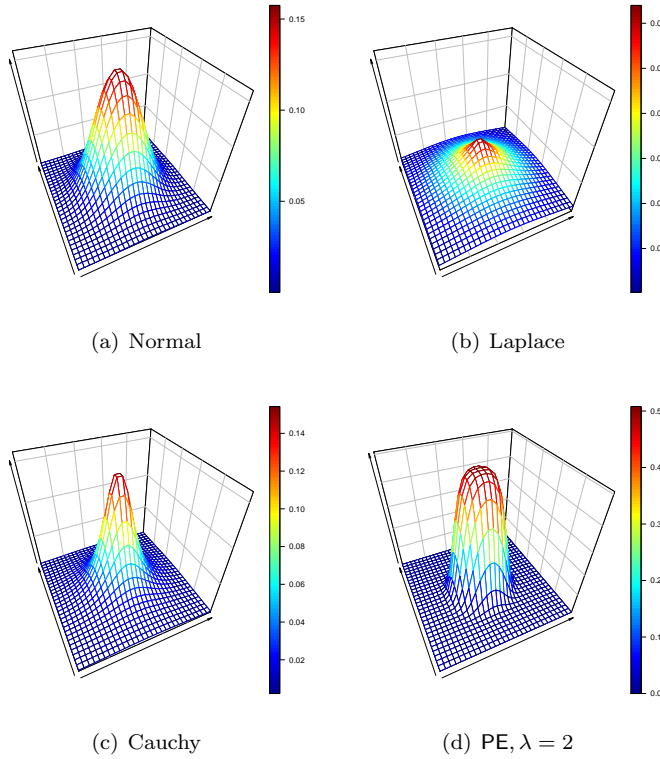


FIGURA 3. Funciones de densidad del vector $\mathbf{X} \sim \text{EC}_2(\mathbf{0}, \mathbf{I}; g)$ para las distribuciones normal, Laplace, Cauchy y exponencial potencia con $\lambda = 2$.

EJEMPLO 1.53 (Distribución t de Student). La función generadora de densidad de un vector aleatorio con distribución t de Student asume la forma

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

Para la distribución t de Student, tenemos que $R^2/p \sim F_{p,\nu}$. Además, la función característica de $\mathbf{X} \sim t_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$ es dada por

$$\varphi(\mathbf{t}) = \frac{\|\sqrt{\nu}\boldsymbol{\Sigma}^{1/2}\mathbf{t}\|^{\nu/2}}{2^{\nu/2-1}\Gamma(\nu/2)} \exp\{i\mathbf{t}^\top \boldsymbol{\mu}\} K_{\nu/2}(\|\sqrt{\nu}\boldsymbol{\Sigma}^{1/2}\mathbf{t}\|), \quad \mathbf{t} \in \mathbb{R}^p,$$

donde $K_\nu(x)$ denota la función de Bessel modificada de segundo tipo. Un caso particular importante corresponde a la distribución Cauchy, cuando $\nu = 1$, mientras que la distribución normal corresponde al caso límite $\nu \rightarrow \infty$.

EJEMPLO 1.54 (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 $R^{2\lambda} \sim \text{Gama}(\frac{1}{2}, \frac{p}{2\lambda})$. Debemos destacar que esta clase de distribuciones contiene la distribución normal como un caso particular cuando $\lambda = 1$. Mientras que tiene colas más pesadas que la normal si $\lambda < 1$ y colas más livianas para el caso $\lambda > 1$. Otro caso particular de interés es la distribución Laplace, que es recuperada cuando $\lambda = 1/2$.

DEFINICIÓN 1.55. Sea $\boldsymbol{\mu} \in \mathbb{R}^p$, $\boldsymbol{\Sigma}$ matriz $p \times p$ definida positiva y \mathbf{H} función de distribución de la 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

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

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

EJEMPLO 1.56 (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 1.57. Un vector aleatorio $\mathbf{X} \sim \text{SMN}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}; \mathbf{H})$ 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 $W \sim \mathbf{H}(\nu)$ son independientes. También podemos utilizar la siguiente estructura jerárquica:

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

Esta representación permite, por ejemplo

$$\mathbf{E}(\mathbf{X}) = \mathbf{E}(\mathbf{E}(\mathbf{X}|W)) = \boldsymbol{\mu}$$

$$\text{Cov}(\mathbf{X}) = \mathbf{E}(\text{Cov}(\mathbf{X}|W)) + \text{Cov}(\mathbf{E}(\mathbf{X}|W)) = \mathbf{E}(W^{-1})\boldsymbol{\Sigma},$$

siempre que $\mathbf{E}(W^{-1}) < \infty$.

EJEMPLO 1.58 (Distribución t de Student). Para $\mathbf{X} \sim t_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$, con $\nu > 0$, podemos escribir

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

es decir, la función de densidad asociado a la variable de mezcla, es dada por

$$h(\omega; \nu) = \frac{(\nu/2)^{\nu/2} \omega^{\nu/2-1}}{\Gamma(\nu/2)} \exp(-\nu\omega/2).$$

EJEMPLO 1.59 (Distribución normal contaminada). Considere $\mathbf{X} \sim \text{CN}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \epsilon, \gamma)$ [Little \(1988\)](#) donde $0 \leq \epsilon \leq 1$ denota el *porcentaje de contaminación* y $0 < \gamma < 1$ corresponde a un *factor de inflación de escala*. En este caso, la variable de mezcla tiene densidad

$$h(\omega; \boldsymbol{\delta}) = \begin{cases} \epsilon, & \omega = \gamma \\ 1 - \epsilon & \omega = 1 \end{cases},$$

con $\boldsymbol{\delta} = (\epsilon, \gamma)^\top$. Podemos notar que la función de densidad adopta la forma:

$$f(\mathbf{x}) = (1 - \epsilon)|2\pi\boldsymbol{\Sigma}|^{-1/2} \exp(-u/2) + \epsilon\gamma^{p/2}|2\pi\boldsymbol{\Sigma}|^{-1/2} \exp(-\lambda u/2).$$

1.7. Algunas distribuciones no centrales

Las distribuciones chi-cuadrado, F , t de Student no central son derivadas desde la distribución normal multivariada y son útiles para desarrollar la inferencia en modelo de regresión lineal.

RESULTADO 1.60. Sea $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$ y sea $U = \mathbf{Z}^\top \mathbf{Z}$. Entonces $U \sim \chi^2(p)$, con función de densidad

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

DEMOSTRACIÓN. Como U es una función de variables aleatorias normales, entonces su función característica asume la forma

$$\begin{aligned} \varphi_U(t) &= \mathbf{E}\{\exp(itU)\} = \int_{\mathbb{R}^p} \exp(it\mathbf{z}^\top \mathbf{z}) (2\pi)^{-p/2} \exp(-\tfrac{1}{2}\mathbf{z}^\top \mathbf{z}) d\mathbf{z} \\ &= (2\pi)^{-p/2} \int_{\mathbb{R}^p} \exp(-\tfrac{1}{2}(1 - 2it)\mathbf{z}^\top \mathbf{z}) d\mathbf{z} = (1 - 2it)^{-p/2}, \end{aligned}$$

que corresponde a la función característica de una variable aleatoria chi-cuadrado con p grados de libertad. \square

DEFINICIÓN 1.61 (Distribución chi-cuadrado no central). Si $\mathbf{Y} \sim \mathbf{N}_p(\boldsymbol{\mu}, \mathbf{I})$, entonces $U = \mathbf{Y}^\top \mathbf{Y}$ tiene *distribución chi-cuadrado no central* con p grados de libertad y parámetro de no centralidad $\lambda = \boldsymbol{\mu}^\top \boldsymbol{\mu}/2$, en cuyo caso anotamos $U \sim \chi^2(p; \lambda)$.

RESULTADO 1.62. Sea $\mathbf{Y} \sim \mathbf{N}_p(\boldsymbol{\mu}, \mathbf{I})$ donde $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p) \neq \mathbf{0}$ y sea $U = \mathbf{Y}^\top \mathbf{Y}$. Entonces la función característica de U es dada por

$$\varphi_U(t) = (1 - 2it)^{-p/2} \exp\left(\frac{2it\lambda}{1 - 2it}\right),$$

con $\lambda = \boldsymbol{\mu}^\top \boldsymbol{\mu}/2$.

DEMOSTRACIÓN. Como Y_1, \dots, Y_p son variables aleatorias independientes, tenemos

$$\begin{aligned} \varphi_U(t) &= \mathbb{E} \left\{ \exp \left(t \sum_{j=1}^n Y_j^2 \right) \right\} = \mathbb{E} \left\{ \prod_{j=1}^p \exp(tY_j^2) \right\} = \prod_{j=1}^p \mathbb{E} \{ \exp(tY_j^2) \} \\ &= \prod_{j=1}^p \varphi_{Y_j^2}(t). \end{aligned}$$

Ahora, la función característica asociada a la variable aleatoria Y_j^2 es dada por

$$\begin{aligned} \varphi_{Y_j^2}(t) &= \int_{-\infty}^{\infty} \exp(it y_j^2) (2\pi)^{-1/2} \exp\left\{-\frac{1}{2}(y_j - \mu_j)^2\right\} dy_j \\ &= \exp\left\{\frac{\mu_j^2}{2} \left(\frac{1}{1-2it}\right) - \frac{\mu_j^2}{2}\right\} \int_{-\infty}^{\infty} (2\pi)^{-1/2} \exp\left\{-\frac{(1-2it)}{2} \left(y_j - \frac{\mu_j}{1-2it}\right)^2\right\} dy_j, \end{aligned}$$

de este modo,

$$\varphi_{Y_j^2}(t) = (1 - 2it)^{-1/2} \exp\left\{\frac{\mu_j^2}{2} \left(\frac{2it}{1 - 2it}\right)\right\},$$

y por tanto la función característica de la variable $U = \sum_{j=1}^p Y_j^2$, asume la forma

$$\varphi_U(t) = (1 - 2it)^{-p/2} \exp\left(\frac{2it\lambda}{1 - 2it}\right), \quad \lambda = \boldsymbol{\mu}^\top \boldsymbol{\mu}/2.$$

□

OBSERVACIÓN 1.63. Es interesante notar que la función característica de la variable $U = \mathbf{Y}^\top \mathbf{Y}$, puede ser escrita como

$$\begin{aligned} \varphi_U(t) &= (1 - 2it)^{-p/2} \exp\left(\frac{\lambda}{1 - 2it} - \lambda\right) \\ &= (1 - 2it)^{-p/2} e^{-\lambda} \sum_{k=0}^{\infty} \frac{\{\lambda/(1 - 2it)\}^k}{k!} \\ &= \sum_{k=0}^{\infty} \frac{e^{-\lambda} \lambda^k}{k!} (1 - 2it)^{-(p+2k)/2}. \end{aligned}$$

Es decir, la función característica de U es un *promedio ponderado con pesos Poisson* de funciones características de variables aleatorias chi-cuadrado con $p + 2k$ grados de libertad.

Usando la relación entre funciones características y sus correspondientes funciones de densidad, sigue que la chi-cuadrado no central tiene la siguiente representación de mezcla

$$U|Z \sim \chi^2(p + 2z), \quad Z \sim \text{Poisson}(\lambda), \quad (1.6)$$

con densidad

$$f(u) = \sum_{k=0}^{\infty} \frac{e^{-\lambda} \lambda^k}{k!} \frac{1}{2^{p/2+k} \Gamma(\frac{p}{2} + k)} u^{p/2+k-1} \exp(-u/2), \quad u > 0.$$

La representación en (1.6) es muy útil para obtener los momentos de una variable aleatoria con distribución chi-cuadrado no central. En efecto, el valor esperado de $U \sim \chi^2(p; \lambda)$ es dado por

$$\mathbb{E}(U) = \mathbb{E}\{\mathbb{E}(U|Z)\} = \mathbb{E}\{p + 2Z\} = p + 2\mathbb{E}(Z) = p + 2\lambda,$$

mientras que la varianza de U puede ser calculada como

$$\begin{aligned} \text{var}(U) &= \mathbb{E}\{\text{var}(U|Z)\} + \text{var}\{\mathbb{E}(U|Z)\} \\ &= \mathbb{E}\{2(p + 2Z)\} + \text{var}(p + 2Z) \\ &= 2p + 4\lambda + 4\lambda = 2p + 8\lambda. \end{aligned}$$

RESULTADO 1.64. Si $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ donde $\boldsymbol{\Sigma}$ es matriz no singular. Entonces

- (a) $(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) \sim \chi^2(p)$.
- (b) $\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} \sim \chi^2(p; \lambda)$, donde $\lambda = \frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$.

DEMOSTRACIÓN. La idea de la demostración se basa en transformar los componentes de \mathbf{X} en variables aleatorias normales independientes. Considere $\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}^\top$ con \mathbf{B} no singular. Para probar (a), tome

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

luego $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$ y de este modo

$$(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) = \mathbf{Z}^\top \mathbf{Z} \sim \chi^2(p; 0).$$

Para probar (b), sea $\mathbf{Y} = \mathbf{B}^{-1} \mathbf{X}$, luego

$$\mathbf{Y} \sim \mathbf{N}_p(\mathbf{B}^{-1} \boldsymbol{\mu}, \mathbf{I}),$$

y

$$\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} = \mathbf{Y}^\top \mathbf{B}^\top \boldsymbol{\Sigma}^{-1} \mathbf{B} \mathbf{Y} = \mathbf{Y}^\top \mathbf{Y},$$

que por definición tiene una distribución chi-cuadrado no central, con parámetro de no centralidad

$$\lambda = \frac{1}{2} (\mathbf{B}^{-1} \boldsymbol{\mu})^\top (\mathbf{B}^{-1} \boldsymbol{\mu}) = \frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}.$$

□

DEFINICIÓN 1.65 (Distribución F no central). Sea $X_1 \sim \chi^2(\nu_1; \lambda)$ y $X_2 \sim \chi^2(\nu_2)$ variables aleatorias independientes. Entonces,

$$F = \frac{X_1/\nu_1}{X_2/\nu_2} \sim F(\nu_1, \nu_2, \lambda),$$

es decir F sigue una *distribución F no central* con ν_1 y ν_2 grados de libertad y parámetro de no centralidad λ .

DEFINICIÓN 1.66 (Distribución Beta no central). Considere $U_1 \sim \chi^2(\nu_1, \lambda)$, $U_2 \sim \chi^2(\nu_2)$ tal que U_1 y U_2 son variables aleatorias independientes. Entonces,

$$G = \frac{U_1}{U_1 + U_2} \sim \text{Beta}(\nu_1, \nu_2, \lambda),$$

esto es, G sigue una *distribución Beta no central* con parámetros de forma y escala ν_1 y ν_2 , respectivamente y parámetro de no centralidad λ .

DEFINICIÓN 1.67 (Distribución t de Student no central). Si $Y \sim \mathbf{N}(\mu, \sigma^2)$ y $U/\sigma^2 \sim \chi^2(\nu)$ son independientes, entonces

$$T = \frac{Y}{\sqrt{U/\nu}} \sim t_\nu(\lambda), \quad \lambda = \mu/\sigma,$$

es llamada una variable aleatoria con *distribución t de Student no central* con ν grados de libertad y parámetro de no centralidad λ .

Note también que si $Z \sim \mathbf{N}(0, 1)$, $U \sim \chi^2(\nu)$, δ es una constante, y Z es independiente de U , entonces

$$T = \frac{Z + \delta}{\sqrt{U/\nu}} \sim t_\nu(\delta).$$

Además el cuadrado de una variable aleatoria t no central se distribuye como una variable aleatoria F no central con parámetro de no centralidad $\delta = \lambda^2/2$. De este modo,

$$t_\nu^2(\lambda) \stackrel{d}{=} F(1, \nu, \lambda^2/2).$$

1.8. Distribución de formas cuadráticas

Para motivar ideas, sabemos que si $\mathbf{Z} \sim \mathbf{N}_p(\mathbf{0}, \mathbf{I})$, entonces $U = \mathbf{Z}^\top \mathbf{Z} \sim \chi^2(p)$ pues corresponde a la suma de variables aleatorias IID $\mathbf{N}(0, 1)$. El objetivo de esta sección es proveer condiciones bajo las cuales variables aleatorias de la forma $U = \mathbf{X}^\top \mathbf{A} \mathbf{X}$ con $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ siguen una distribución chi-cuadrado no central así como establecer la independencia entre dos o más formas cuadráticas.

RESULTADO 1.68. Si $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \mathbf{I})$ y $\mathbf{A} \in \mathbb{R}^{p \times p}$ es matriz simétrica. Entonces $\mathbf{X}^\top \mathbf{A} \mathbf{X} \sim \chi^2(k; \theta)$ si y sólo si \mathbf{A} es idempotente, en cuyo caso los grados de libertad y el parámetro de no centralidad están dados por

$$k = \text{rg}(\mathbf{A}) = \text{tr}(\mathbf{A}), \quad y \quad \theta = \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu},$$

respectivamente.

DEMOSTRACIÓN. Suponga que \mathbf{A} es idempotente de rango k . Entonces existe una matriz ortogonal \mathbf{P} tal que

$$\mathbf{P}^\top \mathbf{A} \mathbf{P} = \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

Sea $\mathbf{Y} = \mathbf{P}^\top \mathbf{X}$, entonces $\mathbf{Y} \sim \mathbf{N}_p(\mathbf{P}^\top \boldsymbol{\mu}, \mathbf{I})$, y

$$\mathbf{X}^\top \mathbf{A} \mathbf{X} = \mathbf{Y}^\top \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Y} = \sum_{i=1}^k Y_i^2,$$

que sigue una distribución chi-cuadrado con k grados de libertad. Para el parámetro de no centralidad θ , note que

$$\begin{aligned} \mathbb{E}\{\chi^2(k; \theta)\} &= k + 2\theta = \mathbb{E}(\mathbf{X}^\top \mathbf{A} \mathbf{X}) = \text{tr}(\mathbb{E}(\mathbf{X} \mathbf{X}^\top) \mathbf{A}) \\ &= \text{tr}((\mathbf{I} + \boldsymbol{\mu} \boldsymbol{\mu}^\top) \mathbf{A}) = k + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}, \end{aligned}$$

y de ahí que $\theta = \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}$.

Ahora, suponga que $\mathbf{X}^\top \mathbf{A} \mathbf{X} \sim \chi^2(k; \theta)$. Si \mathbf{A} tiene rango r , entonces para \mathbf{P} matriz ortogonal $p \times p$,

$$\mathbf{P}^\top \mathbf{A} \mathbf{P} = \begin{pmatrix} \mathbf{\Lambda}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix},$$

con $\mathbf{\Lambda}_1 = \text{diag}(\lambda_1, \dots, \lambda_r)$, donde $\lambda_1, \dots, \lambda_r$ son los valores propios no nulos de \mathbf{A} . Sea $\mathbf{Y} = \mathbf{P}^\top \mathbf{X}$, entonces

$$\mathbf{X}^\top \mathbf{A} \mathbf{X} = \mathbf{Y}^\top \mathbf{P}^\top \mathbf{A} \mathbf{P} \mathbf{Y} = \sum_{j=1}^r \lambda_j Y_j^2 = U.$$

Tenemos que $\mathbf{Y} \sim \mathbf{N}_p(\boldsymbol{\delta}, \mathbf{I})$ con $\boldsymbol{\delta} = \mathbf{P}^\top \boldsymbol{\mu}$, de modo que $Y_j^2 \sim \chi^2(1; \delta_j^2/2)$ con función característica

$$\varphi_{Y_j^2}(t) = (1 - 2it)^{-1/2} \exp\left(\frac{it\delta_j^2}{1 - 2it}\right),$$

por la independencia de Y_1, \dots, Y_r sigue que

$$\begin{aligned} \varphi_U(t) &= \prod_{j=1}^r (1 - 2it\lambda_j)^{-1/2} \exp\left(\frac{it\lambda_j\delta_j^2}{1 - 2it\lambda_j}\right) \\ &= \exp\left(it \sum_{j=1}^r \frac{\lambda_j\delta_j^2}{1 - 2it\lambda_j}\right) \prod_{j=1}^r (1 - 2it\lambda_j)^{-1/2}. \end{aligned}$$

Como $\mathbf{X}^\top \mathbf{A} \mathbf{X} \sim \chi_k^2(\theta)$ tiene función característica

$$\varphi_{\mathbf{X}^\top \mathbf{A} \mathbf{X}}(t) = (1 - 2it)^{-k/2} \exp\left(\frac{2it\theta}{1 - 2it}\right),$$

entonces desde las dos expresiones anteriores debemos tener $r = k$, $\lambda_j = 1$, $\forall j$ y $\theta = \sum_j \delta_j^2/2$. Consecuentemente $\mathbf{P}^\top \mathbf{A} \mathbf{P}$ tiene la forma

$$\mathbf{P}^\top \mathbf{A} \mathbf{P} = \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix},$$

que es idempotente. Luego

$$\mathbf{P}^\top \mathbf{A} \mathbf{P} = (\mathbf{P}^\top \mathbf{A} \mathbf{P})(\mathbf{P}^\top \mathbf{A} \mathbf{P}) = \mathbf{P}^\top \mathbf{A}^2 \mathbf{P} \implies \mathbf{A}^2 = \mathbf{A}.$$

□

RESULTADO 1.69. Si $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ donde $\boldsymbol{\Sigma}$ es no singular y \mathbf{X} , $\boldsymbol{\mu}$ y $\boldsymbol{\Sigma}$ son particionados como

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix},$$

donde \mathbf{X}_1 , $\boldsymbol{\mu}_1$ son $k \times 1$ y $\boldsymbol{\Sigma}_{11}$ es $k \times k$. Entonces

$$U = (\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}) - (\mathbf{X}_1 - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{X}_1 - \boldsymbol{\mu}_1) \sim \chi^2(p - k).$$

DEMOSTRACIÓN. Considere $\boldsymbol{\Sigma} = \mathbf{B} \mathbf{B}^\top$, donde \mathbf{B} es no singular y particione \mathbf{B} como

$$\mathbf{B} = \begin{pmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{pmatrix}, \quad \mathbf{B}_1 \in \mathbb{R}^{k \times p}.$$

Luego,

$$\boldsymbol{\Sigma} = \mathbf{B} \mathbf{B}^\top = \begin{pmatrix} \mathbf{B}_1 \mathbf{B}_1^\top & \mathbf{B}_1 \mathbf{B}_2^\top \\ \mathbf{B}_2 \mathbf{B}_1^\top & \mathbf{B}_2 \mathbf{B}_2^\top \end{pmatrix},$$

de donde sigue que $\Sigma_{11} = B_1 B_1^\top$. Ahora, sea $Z = B^{-1}(X - \mu) \sim N_p(0, I)$. De este modo,

$$\begin{pmatrix} B_1 \\ B_2 \end{pmatrix} Z = \begin{pmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \end{pmatrix}.$$

Entonces

$$\begin{aligned} U &= Z^\top Z - Z^\top B_1^\top (B_1 B_1^\top)^{-1} B_1 Z = Z^\top (I - B_1^\top (B_1 B_1^\top)^{-1} B_1) Z \\ &= Z^\top (I - H_1) Z, \quad \text{con } H_1 = B_1^\top (B_1 B_1^\top)^{-1} B_1. \end{aligned}$$

Note que H_1 es simétrica e idempotente y por tanto también lo es $C = I - H_1$. De donde sigue que $U \sim \chi^2(\nu)$, con $\nu = \text{rg}(C) = p - k$. \square

El Resultado 1.68 se puede generalizar al caso que X tiene una matriz de covarianza arbitraria. Suponga que $X \sim N_p(0, \Sigma)$. Una condición para que $X^\top A X$ tenga una distribución chi-cuadrado es

$$\Sigma A \Sigma A = \Sigma A,$$

en cuyo caso los grados de libertad son $k = \text{rg}(A \Sigma)$. Si Σ es no singular, la condición resulta $A \Sigma A = A$.

RESULTADO 1.70. Si $X \sim N_p(0, \Sigma)$ donde Σ tiene rango k ($\leq p$) y si A es una inversa generalizada de Σ ($\Sigma A \Sigma = \Sigma$), entonces $X^\top A X \sim \chi^2(k)$.

DEMOSTRACIÓN. Considere $Y = B X$ donde B es una matriz no singular $p \times p$ tal que

$$B \Sigma B^\top = \begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix}.$$

Particionando $Y = (Y_1^\top, Y_2^\top)^\top$ donde Y_1 es un vector $k \times 1$ sigue que $Y_1 \sim N_k(0, I)$ y $Y_2 = 0$ con probabilidad 1. Es decir, tenemos que

$$Y = (Y_1^\top, 0)^\top, \quad \text{con probabilidad 1.}$$

Ahora, note que

$$\begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix} = B \Sigma B^\top = B \Sigma A \Sigma B^\top$$

pues A es una inversa generalizada de Σ . De este modo,

$$\begin{aligned} \begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix} &= B \Sigma B^\top B^{-\top} A B^{-1} B \Sigma B^\top \\ &= \begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix} B^{-\top} A B^{-1} \begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix}. \end{aligned}$$

Luego, con probabilidad uno,

$$\begin{aligned} X^\top A X &= Y^\top B^{-\top} A B^{-1} Y = (Y_1^\top, 0) B^{-\top} A B^{-1} \begin{pmatrix} Y_1 \\ 0 \end{pmatrix} \\ &= (Y_1^\top, 0) \begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix} B^{-\top} A B^{-1} \begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} Y_1 \\ 0 \end{pmatrix} \\ &= (Y_1^\top, 0) \begin{pmatrix} I_k & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} Y_1 \\ 0 \end{pmatrix} = Y_1^\top Y_1 \sim \chi^2(k). \end{aligned}$$

\square

RESULTADO 1.71. Si $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, donde $\boldsymbol{\Sigma}$ es no singular, y \mathbf{A} es una matriz simétrica $p \times p$. Entonces $\mathbf{X}^\top \mathbf{A} \mathbf{X} \sim \chi^2(k; \lambda)$, donde $k = \text{rg}(\mathbf{A})$, $\lambda = \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} / 2$ si y sólo si $\mathbf{A} \boldsymbol{\Sigma}$ es matriz idempotente.

DEMOSTRACIÓN. Considere $\mathbf{Y} = \mathbf{B} \mathbf{X}$, donde \mathbf{B} es una matriz no singular $p \times p$ tal que $\mathbf{B} \boldsymbol{\Sigma} \mathbf{B}^\top = \mathbf{I}_p$. Entonces

$$\mathbf{X}^\top \mathbf{A} \mathbf{X} = \mathbf{Y}^\top \mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1} \mathbf{Y},$$

donde $\mathbf{Y} \sim \mathbf{N}_p(\mathbf{B} \boldsymbol{\mu}, \mathbf{I})$. Desde el Resultado 1.68 sigue que $\mathbf{X}^\top \mathbf{A} \mathbf{X}$ tiene distribución chi-cuadrado sólo si $\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1}$ es idempotente. Esto es equivalente a mostrar que $\mathbf{A} \boldsymbol{\Sigma}$ es idempotente.

Si $\mathbf{A} \boldsymbol{\Sigma}$ es idempotente, tenemos

$$\mathbf{A} = \mathbf{A} \boldsymbol{\Sigma} \mathbf{A} = \mathbf{A} \mathbf{B}^{-1} \mathbf{B}^{-\top} \mathbf{A}, \quad (\boldsymbol{\Sigma} = \mathbf{B}^{-1} \mathbf{B}^{-\top})$$

así, pre- y post-multiplicando por $\mathbf{B}^{-\top}$ y \mathbf{B}^{-1} , obtenemos

$$\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1} = (\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1})(\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1}),$$

y por tanto es idempotente.

Por otro lado, si $\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1}$ es idempotente, entonces

$$\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1} = (\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1})(\mathbf{B}^{-\top} \mathbf{A} \mathbf{B}^{-1}) = \mathbf{B}^{-\top} \mathbf{A} \boldsymbol{\Sigma} \mathbf{A} \mathbf{B}^{-1},$$

es decir $\mathbf{A} = \mathbf{A} \boldsymbol{\Sigma} \mathbf{A}$ y de ahí que $\mathbf{A} \boldsymbol{\Sigma}$ es idempotente. \square

EJEMPLO 1.72. Sea X_1, \dots, X_n variables aleatorias IID $\mathbf{N}(\theta, \sigma^2)$, en este caso podemos definir $\mathbf{X} = (X_1, \dots, X_n)^\top$ tal que $\mathbf{X} \sim \mathbf{N}_n(\theta \mathbf{1}_n, \sigma^2 \mathbf{I})$. Considere la forma cuadrática

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

con $\mathbf{C} = \mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}^\top$ y $\mathbf{A} = \mathbf{C} / \sigma^2$. De esta manera

$$\mathbf{A} \boldsymbol{\Sigma} = \mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}^\top,$$

que es idempotente. Además

$$\text{rg}(\mathbf{A}) = \text{tr}(\mathbf{C}) = \text{tr}\left(\mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}^\top\right) = n - 1,$$

y

$$\lambda = \frac{\theta^2}{2} \mathbf{1}^\top \mathbf{A} \mathbf{1} = \frac{\theta^2}{2\sigma^2} \mathbf{1}^\top \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}^\top\right) \mathbf{1} = 0.$$

Finalmente,

$$Q = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi^2(n - 1).$$

RESULTADO 1.73. Sea $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $Q_1 = \mathbf{X}^\top \mathbf{A} \mathbf{X}$ y $Q_2 = \mathbf{X}^\top \mathbf{B} \mathbf{X}$. Entonces Q_1 y Q_2 son independientes si y sólo si $\mathbf{A} \boldsymbol{\Sigma} \mathbf{B} = \mathbf{0}$.

DEMOSTRACIÓN. Tenemos $\boldsymbol{\Sigma} = \mathbf{T} \mathbf{T}^\top$, y defina $\mathbf{G}_1 = \mathbf{T}^\top \mathbf{A} \mathbf{T}$, $\mathbf{G}_2 = \mathbf{T}^\top \mathbf{B} \mathbf{T}$. Note que si $\mathbf{A} \boldsymbol{\Sigma} \mathbf{B} = \mathbf{0}$, entonces

$$\mathbf{G}_1 \mathbf{G}_2 = (\mathbf{T}^\top \mathbf{A} \mathbf{T})(\mathbf{T}^\top \mathbf{B} \mathbf{T}) = \mathbf{T}^\top \mathbf{A} \boldsymbol{\Sigma} \mathbf{B} \mathbf{T} = \mathbf{0}.$$

Debido a la simetría de \mathbf{G}_1 y \mathbf{G}_2 , sigue que

$$\mathbf{0} = (\mathbf{G}_1 \mathbf{G}_2)^\top = \mathbf{G}_2^\top \mathbf{G}_1^\top = \mathbf{G}_2 \mathbf{G}_1.$$

Como $\mathbf{G}_1 \mathbf{G}_2 = \mathbf{G}_2 \mathbf{G}_1$ existe una matriz ortogonal \mathbf{P} que simultáneamente diagonaliza \mathbf{G}_1 y \mathbf{G}_2 , esto es:

$$\mathbf{P}^\top \mathbf{G}_1 \mathbf{P} = \mathbf{P}^\top \mathbf{T}^\top \mathbf{A} \mathbf{P} = \mathbf{D}_1,$$

$$\mathbf{P}^\top \mathbf{G}_2 \mathbf{P} = \mathbf{P}^\top \mathbf{T}^\top \mathbf{B} \mathbf{P} = \mathbf{D}_2.$$

De este modo,

$$\mathbf{0} = \mathbf{G}_1 \mathbf{G}_2 = \mathbf{P} \mathbf{D}_1 \mathbf{P}^\top \mathbf{P} \mathbf{D}_2 \mathbf{P}^\top = \mathbf{P} \mathbf{D}_1 \mathbf{D}_2 \mathbf{P}^\top$$

lo que es verdad si $\mathbf{D}_1 \mathbf{D}_2 = \mathbf{0}$. Como \mathbf{D}_1 y \mathbf{D}_2 son diagonales, sus elementos diagonales deben ocurrir en posiciones diferentes. Es decir, podemos escribir

$$\mathbf{D}_1 = \begin{pmatrix} M_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \mathbf{D}_2 = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & M_2 \end{pmatrix}.$$

Sea $\mathbf{Y} = \mathbf{P}^\top \mathbf{T}^{-1} \mathbf{X}$, entonces

$$Q_1 = \mathbf{X}^\top \mathbf{A} \mathbf{X} = \mathbf{X}^\top \mathbf{T}^{-\top} \mathbf{P} \mathbf{P}^\top \mathbf{T}^\top \mathbf{A} \mathbf{P} \mathbf{P}^\top \mathbf{T}^{-1} \mathbf{X} = \mathbf{Y}^\top \mathbf{D}_1 \mathbf{Y},$$

$$Q_2 = \mathbf{X}^\top \mathbf{B} \mathbf{X} = \mathbf{X}^\top \mathbf{T}^{-\top} \mathbf{P} \mathbf{P}^\top \mathbf{T}^\top \mathbf{B} \mathbf{P} \mathbf{P}^\top \mathbf{T}^{-1} \mathbf{X} = \mathbf{Y}^\top \mathbf{D}_2 \mathbf{Y}.$$

Además,

$$\text{Cov}(\mathbf{Y}) = \text{Cov}(\mathbf{P}^\top \mathbf{T}^{-1} \mathbf{X}) = \mathbf{P}^\top \mathbf{T}^{-1} \text{Cov}(\mathbf{X}) \mathbf{T}^{-\top} \mathbf{P} = \mathbf{I}.$$

En efecto, $\mathbf{Y} \sim \mathbf{N}_p(\mathbf{P}^\top \mathbf{T}^{-1} \boldsymbol{\mu}, \mathbf{I})$. Ahora, particionando adecuadamente \mathbf{Y} , sigue que

$$\mathbf{Y}^\top \mathbf{D}_1 \mathbf{Y} = (\mathbf{Y}_1^\top, \mathbf{Y}_2^\top) \begin{pmatrix} M_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \end{pmatrix} = \mathbf{Y}_1^\top M_1 \mathbf{Y}_1,$$

$$\mathbf{Y}^\top \mathbf{D}_2 \mathbf{Y} = (\mathbf{Y}_1^\top, \mathbf{Y}_2^\top) \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & M_2 \end{pmatrix} \begin{pmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \end{pmatrix} = \mathbf{Y}_2^\top M_2 \mathbf{Y}_2,$$

y la independencia entre Q_1 y Q_2 sigue desde la independencia entre \mathbf{Y}_1 y \mathbf{Y}_2 . \square

RESULTADO 1.74. Sea $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $Q = \mathbf{X}^\top \mathbf{A} \mathbf{X}$ y $U = \mathbf{B} \mathbf{X}$. Entonces Q y U son independientes si y sólo si $\mathbf{B} \boldsymbol{\Sigma} \mathbf{A} = \mathbf{0}$.

EJEMPLO 1.75. Considere X_1, \dots, X_n muestra aleatoria desde $\mathbf{N}(\theta, \sigma^2)$, así

$$\mathbf{X} = (X_1, \dots, X_n)^\top \sim \mathbf{N}_n(\theta \mathbf{1}, \sigma^2 \mathbf{I}_n).$$

Tenemos

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

Como $\mathbf{C} \mathbf{1} = \mathbf{0}$ sigue la independencia entre \bar{X} y S^2 .

Considere los siguientes dos lemas, los que permitirán mostrar el resultado principal de esta sección.

LEMA 1.76. Sean $\mathbf{A}_1, \dots, \mathbf{A}_k$ matrices $m \times m$ simétricas e idempotentes y suponga que

$$\mathbf{A}_1 + \mathbf{A}_2 + \dots + \mathbf{A}_k = \mathbf{I}_m.$$

Entonces $\mathbf{A}_i \mathbf{A}_j = \mathbf{0}$ para todo $i \neq j$.

DEMOSTRACIÓN. Considere cualquiera de esas matrices, digamos \mathbf{A}_h y denote su rango por r . Como \mathbf{A}_h es simétrica e idempotente, existe una matriz ortogonal \mathbf{P} tal que

$$\mathbf{P}^\top \mathbf{A}_h \mathbf{P} = \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

Para $j \neq h$, defina $\mathbf{B}_j = \mathbf{P}^\top \mathbf{A}_j \mathbf{P}$, y note que

$$\mathbf{I}_m = \mathbf{P}^\top \mathbf{P} = \mathbf{P}^\top \left(\sum_{j=1}^k \mathbf{A}_j \right) \mathbf{P} = \sum_{j=1}^k \mathbf{P}^\top \mathbf{A}_j \mathbf{P} = \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + \sum_{j \neq h} \mathbf{B}_j.$$

O equivalentemente,

$$\sum_{j \neq h} \mathbf{B}_j = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{m-r} \end{pmatrix}.$$

Claramente, dado que \mathbf{A}_j es simétrica e idempotente, sigue que \mathbf{B}_j también lo es. De modo que, sus elementos diagonales son no negativos. Además, $(\mathbf{B}_j)_{ll} = 0$, para $l = 1, \dots, r$. Así, sigue que \mathbf{B}_j debe ser de la forma

$$\mathbf{B}_j = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_j \end{pmatrix},$$

donde \mathbf{C}_j es matriz $(m-r) \times (m-r)$, simétrica e idempotente. Ahora, para cualquier $j \neq h$

$$\mathbf{P}^\top \mathbf{A}_h \mathbf{A}_j \mathbf{P} = (\mathbf{P}^\top \mathbf{A}_h \mathbf{P})(\mathbf{P}^\top \mathbf{A}_j \mathbf{P}) = \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_j \end{pmatrix} = \mathbf{0},$$

lo que es verdad, sólo si $\mathbf{A}_h \mathbf{A}_j = \mathbf{0}$, pues \mathbf{P} es no singular. Notando que h es arbitrario, la prueba es completa. \square

LEMA 1.77. Sean $\mathbf{A}_1, \dots, \mathbf{A}_k$ matrices simétricas de orden $m \times m$ y defina

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2 + \dots + \mathbf{A}_k.$$

Considere las siguientes afirmaciones,

- (a) \mathbf{A}_i es idempotente, para $i = 1, \dots, k$.
- (b) \mathbf{A} es idempotente.
- (c) $\mathbf{A}_i \mathbf{A}_j = \mathbf{0}$, para $i \neq j$.

Entonces, si dos condiciones son satisfechas, la tercera condición debe ser verdadera.

DEMOSTRACIÓN. Primero mostraremos que (a) y (b) implica (c). Como \mathbf{A} es simétrica e idempotente, existe una matriz ortogonal \mathbf{P} tal que

$$\mathbf{P}^\top \mathbf{A} \mathbf{P} = \mathbf{P}^\top (\mathbf{A}_1 + \dots + \mathbf{A}_k) \mathbf{P} = \begin{pmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad (1.7)$$

donde $r = \text{rg}(\mathbf{A})$.

Sea $\mathbf{B}_i = \mathbf{P}^\top \mathbf{A}_i \mathbf{P}$, para $i = 1, \dots, k$, y note que \mathbf{B}_i es simétrica e idempotente. Es decir, \mathbf{B}_i debe ser de la forma

$$\mathbf{B}_i = \begin{pmatrix} \mathbf{C}_i & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix},$$

donde la matriz $r \times r$, \mathbf{C}_i debe ser simétrica e idempotente. Por (1.7), tenemos

$$\mathbf{C}_1 + \dots + \mathbf{C}_k = \mathbf{I}_r.$$

Por el Lema 1.76, sigue que $\mathbf{C}_i \mathbf{C}_j = \mathbf{0}$ para $i \neq j$, de donde obtenemos $\mathbf{B}_i \mathbf{B}_j = \mathbf{0}$ y de ahí que $\mathbf{A}_i \mathbf{A}_j = \mathbf{0}$, para $i \neq j$.

Que (a) y (c) implican (b), sigue de notar

$$\begin{aligned} \mathbf{A}^2 &= \left(\sum_{i=1}^k \mathbf{A}_i \right)^2 = \sum_{i=1}^k \sum_{j=1}^k \mathbf{A}_i \mathbf{A}_j = \sum_{i=1}^k \mathbf{A}_i^2 + \sum_{i \neq j} \sum \mathbf{A}_i \mathbf{A}_j \\ &= \sum_{i=1}^k \mathbf{A}_i = \mathbf{A}. \end{aligned}$$

Finalmente, para probar que (b) y (c) implican (a). Suponga que (c) es verdad, entonces $\mathbf{A}_i \mathbf{A}_j = \mathbf{A}_j \mathbf{A}_i$ para todo $i \neq j$ y las matrices $\mathbf{A}_1, \dots, \mathbf{A}_k$ pueden ser diagonalizadas simultáneamente. Esto es, existe una matriz ortogonal \mathbf{Q} tal que

$$\mathbf{Q}^\top \mathbf{A}_i \mathbf{Q} = \mathbf{D}_i, \quad i = 1, \dots, k,$$

donde cada una de las matrices $\mathbf{D}_1, \dots, \mathbf{D}_k$ es diagonal. Además,

$$\mathbf{D}_i \mathbf{D}_j = \mathbf{Q}^\top \mathbf{A}_i \mathbf{Q} \mathbf{Q}^\top \mathbf{A}_j \mathbf{Q} = \mathbf{Q}^\top \mathbf{A}_i \mathbf{A}_j \mathbf{Q} = \mathbf{0}, \quad i \neq j. \quad (1.8)$$

Como \mathbf{A} es simétrica e idempotente, también lo es la matriz diagonal

$$\mathbf{Q}^\top \mathbf{A} \mathbf{Q} = \mathbf{D}_1 + \dots + \mathbf{D}_k,$$

y cada elemento diagonal de $\mathbf{Q}^\top \mathbf{A} \mathbf{Q}$ debe ser 0 o 1, y por (1.8), lo mismo es válido para los elementos diagonales de $\mathbf{D}_1, \dots, \mathbf{D}_k$.

De este modo, \mathbf{D}_i es simétrica e idempotente y de ahí que también lo es

$$\mathbf{A}_i = \mathbf{Q} \mathbf{D}_i \mathbf{Q}^\top, \quad i = 1, \dots, k,$$

lo que termina la prueba. \square

OBSERVACIÓN 1.78. Suponga que las condiciones del Lema 1.77 son satisfechas. Entonces (a) implica que $\text{rg}(\mathbf{A}_i) = \text{tr}(\mathbf{A}_i)$, y desde (b), sigue que

$$\text{rg}(\mathbf{A}) = \text{tr}(\mathbf{A}) = \text{tr} \left(\sum_{i=1}^k \mathbf{A}_i \right) = \sum_{i=1}^k \text{tr}(\mathbf{A}_i) = \sum_{i=1}^k \text{rg}(\mathbf{A}_i).$$

RESULTADO 1.79 (Teorema de Cochran). Sea $\mathbf{X} \sim \mathbf{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, con $\boldsymbol{\Sigma} > \mathbf{0}$. Suponga que \mathbf{A}_i , es una matriz simétrica de orden $p \times p$ con rango r_i , para $i = 1, \dots, k$, y

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2 + \dots + \mathbf{A}_k,$$

es de rango r . Considere las condiciones:

- (a) $\mathbf{A}_i \boldsymbol{\Sigma}$ es idempotente, para $i = 1, \dots, k$.
- (b) $\mathbf{A} \boldsymbol{\Sigma}$ es idempotente.
- (c) $\mathbf{A}_i \boldsymbol{\Sigma} \mathbf{A}_j = \mathbf{0}$, para $i \neq j$.
- (d) $r = \sum_{i=1}^k r_i$.

si dos de (a), (b) y (c) se satisfacen, o si (b) y (d) son satisfechas. Entonces,

- (i) $\mathbf{X}^\top \mathbf{A}_i \mathbf{X} \sim \chi^2(r_i; \lambda_i)$, con $\lambda_i = \boldsymbol{\mu}^\top \mathbf{A}_i \boldsymbol{\mu} / 2$, para $i = 1, \dots, k$.
- (ii) $\mathbf{X}^\top \mathbf{A} \mathbf{X} \sim \chi^2(r; \lambda)$, con $\lambda = \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} / 2$.
- (iii) $\mathbf{X}^\top \mathbf{A}_1 \mathbf{X}, \mathbf{X}^\top \mathbf{A}_2 \mathbf{X}, \dots, \mathbf{X}^\top \mathbf{A}_k \mathbf{X}$ son mutuamente independientes.

DEMOSTRACIÓN. Tenemos que $\Sigma = \mathbf{T}\mathbf{T}^\top$ y las condiciones (a)-(d), pueden ser expresadas como:

- (a) $\mathbf{T}^\top \mathbf{A}_i \mathbf{T}$ es idempotente, para $i = 1, \dots, k$.
- (b) $\mathbf{T}^\top \mathbf{A} \mathbf{T}$ es idempotente.
- (c) $(\mathbf{T}^\top \mathbf{A}_i \mathbf{T})(\mathbf{T}^\top \Sigma \mathbf{A}_j \mathbf{T}) = \mathbf{0}$, para $i \neq j$.
- (d) $\text{rg}(\mathbf{T}^\top \mathbf{A} \mathbf{T}) = \sum_{i=1}^k \text{rg}(\mathbf{T}^\top \mathbf{A}_i \mathbf{T})$.

Como $\mathbf{T}^\top \mathbf{A}_1 \mathbf{T}, \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}, \dots, \mathbf{T}^\top \mathbf{A}_k \mathbf{T}$ y $\mathbf{T}^\top \mathbf{A} \mathbf{T}$ satisfacen las condiciones del Lema 1.77 (y de la Observación 1.78). Entonces, las condiciones (a)-(d) se satisfacen.

Sabemos que (a) implica (i) y (b) implica (ii). Mientras que, Resultado 1.73 con (c), garantiza (iii), lo que completa la prueba. \square

Ejercicios

- 1.1 Sean $\mathbf{X}_1, \dots, \mathbf{X}_n$ vectores aleatorios independientes con $\mathbf{X}_i \sim \mathbf{N}_p(\boldsymbol{\mu}, \Sigma)$, para $i = 1, \dots, n$. Obtenga la distribución de

$$\sum_{i=1}^n \alpha_i \mathbf{X}_i,$$

con $\alpha_1, \dots, \alpha_n$ constantes fijas.

- 1.2 Si $\mathbf{X}_1, \dots, \mathbf{X}_n$ son independientes cada uno con $\mathbf{X}_i \sim \mathbf{N}_p(\boldsymbol{\mu}, \Sigma)$. Muestre que la distribución del vector de medias

$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i,$$

es $\mathbf{N}_p(\boldsymbol{\mu}, \frac{1}{n} \Sigma)$.

- 1.3 Demuestre el Resultado 1.20, usando la función característica de un vector aleatorio normal.
- 1.4 Sean X_1, \dots, X_n variables aleatorias independientes e idénticamente distribuidas $\mathbf{N}(\mu, \sigma^2)$ y defina

$$Q = \frac{1}{2(n-1)} \sum_{i=1}^{n-1} (X_{i+1} - X_i)^2,$$

¿Es Q un estimador insesgado de σ^2 ?

- 1.5 Sea $\mathbf{X} \sim \mathbf{N}_n(\boldsymbol{\mu}, \Sigma)$ y defina

$$\mathbf{Y} = \mathbf{T}^\top \Sigma^{-1/2} (\mathbf{X} - \boldsymbol{\mu}), \quad \mathbf{u} = \mathbf{T}^\top \Sigma^{1/2} \mathbf{A} \boldsymbol{\mu}.$$

con \mathbf{T} ortogonal y $\mathbf{A} = \mathbf{A}^\top$. Obtenga la distribución de \mathbf{Y} y calcule $\text{var}(\mathbf{u}^\top \mathbf{Y})$.

- 1.6 Considere \mathbf{Z} matriz aleatoria $n \times p$ con función característica

$$\varphi_{\mathbf{Z}}(\mathbf{T}) = \mathbb{E}\{\exp(i \operatorname{tr}(\mathbf{T}^\top \mathbf{Z}))\} = \exp\{-\frac{1}{2} \operatorname{tr}(\mathbf{T}^\top \mathbf{T})\}.$$

con $\mathbf{T} \in \mathbb{R}^{n \times p}$. Obtenga la función característica de

$$\mathbf{Y} = \boldsymbol{\Sigma}^{1/2} \mathbf{Z} \boldsymbol{\Theta}^{1/2} + \boldsymbol{\mu},$$

donde $\boldsymbol{\mu} \in \mathbb{R}^{n \times p}$ y $\boldsymbol{\Sigma}$, $\boldsymbol{\Theta}$ son matrices semidefinidas positivas $n \times n$ y $p \times p$, respectivamente.

- 1.7 Sea $\mathbf{Z} = \mathbf{U} \mathbf{D} \boldsymbol{\alpha} + \boldsymbol{\epsilon}$ con $\mathbf{U} \in \mathbb{R}^{n \times p}$ tal que $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$, \mathbf{D} es matriz diagonal $p \times p$ y $\boldsymbol{\epsilon} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$. Considere

$$\hat{\boldsymbol{\alpha}} = (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^\top \mathbf{Z}.$$

donde λ es un escalar positivo.

- (a) Obtenga la distribución de $\hat{\boldsymbol{\alpha}}$,
 (b) Muestre que

$$\boldsymbol{\alpha} - \mathbb{E}(\hat{\boldsymbol{\alpha}}) = \lambda(\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \boldsymbol{\alpha}.$$

- 1.8 Sea $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ y considere $\mathbf{b} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$, $\mathbf{u} = (\mathbf{D}^{-1} + \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top (\mathbf{Y} - \mathbf{X}\mathbf{b})$, donde $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Z} \in \mathbb{R}^{n \times q}$ y \mathbf{D} es matriz no singular $q \times q$.

- (a) Halle la distribución de \mathbf{b} y \mathbf{u} ,
 (b) ¿Son \mathbf{b} y \mathbf{u} independientes?

- 1.9 Considere

$$\begin{pmatrix} \mathbf{Y} \\ \mathbf{b} \end{pmatrix} \sim \mathbf{N}_{n+q} \left(\begin{pmatrix} \mathbf{X}\boldsymbol{\beta} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{Z}\mathbf{D}\mathbf{Z}^\top + \mathbf{R} & \mathbf{Z}\mathbf{D} \\ \mathbf{D}\mathbf{Z}^\top & \mathbf{D} \end{pmatrix} \right),$$

donde $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Z} \in \mathbb{R}^{n \times q}$ y \mathbf{R} , \mathbf{D} son matrices no singulares $n \times n$ y $q \times q$, respectivamente. Determine la distribución de $\mathbf{b}|\mathbf{Y}$.

- 1.10 Sea $U_i \sim \chi^2(n_i; \lambda_i)$, $i = 1, \dots, K$ variables aleatorias independientes. Muestre que

$$U = \sum_{i=1}^K U_i \sim \chi^2(n; \lambda),$$

donde $n = \sum_{i=1}^K n_i$ y $\lambda = \sum_{i=1}^K \lambda_i$.

- 1.11 Sea $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$, donde \mathbf{X} es matriz $n \times p$ con $\operatorname{rg}(\mathbf{X}) = p$ y $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$. Defina

$$Q = \frac{(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top [\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top]^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})}{\sigma^2},$$

donde $\mathbf{G} \in \mathbb{R}^{m \times p}$ con $\operatorname{rg}(\mathbf{G}) = m$ y \mathbf{g} es vector m -dimensional. Determine la distribución de Q .

- 1.12 Sea $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$ y considere las formas cuadráticas

$$Q_1 = \frac{\hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}}}{\sigma^2}, \quad Q_2 = \frac{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{\sigma^2},$$

donde $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ con $\mathbf{X} \in \mathbb{R}^{n \times p}$ y $\operatorname{rg}(\mathbf{X}) = p$.

- (a) Halle la distribución de Q_i , $i = 1, 2$.
 (b) Sea $Q = Q_1 + Q_2$, mostrar la independencia conjunta de Q_1 y Q_2 .

- 1.13 Considere $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ con $\mathbf{X} \in \mathbb{R}^{n \times p}$ y $\boldsymbol{\beta} \in \mathbb{R}^p$ y sea $Q = Q_1 + Q_2$, donde

$$Q_1 = \frac{\mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}}{\sigma^2}, \quad Q_2 = \frac{\mathbf{Y}^\top (\mathbf{H} - \frac{1}{n} \mathbf{J}) \mathbf{Y}}{\sigma^2},$$

con $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$. Muestre que Q_1 y Q_2 tienen distribuciones chi-cuadrado independientes.

- 1.14 Considere

$$\mathbf{Y} = (\mathbf{I}_p \otimes \mathbf{1}_n) \boldsymbol{\alpha} + \boldsymbol{\epsilon},$$

donde $\mathbf{Y} = (\mathbf{Y}_1^\top, \mathbf{Y}_2^\top, \dots, \mathbf{Y}_p^\top)^\top$ con \mathbf{Y}_i vector n -dimensional, para $i = 1, \dots, p$, $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_p)^\top$ y $\boldsymbol{\epsilon} \sim \mathbf{N}_{np}(\mathbf{0}, \sigma^2 \mathbf{I}_{np})$. Sean

$$Q_1 = \frac{\mathbf{Y}^\top (\mathbf{I}_p \otimes \frac{1}{n} \mathbf{J}_n) \mathbf{Y}}{\sigma^2}, \quad \text{y} \quad Q_2 = \frac{\mathbf{Y}^\top (\mathbf{I}_p \otimes \mathbf{C}) \mathbf{Y}}{\sigma^2},$$

donde $\mathbf{J}_n = \mathbf{1}_n \mathbf{1}_n^\top$ y $\mathbf{C} = \mathbf{I}_n - \frac{1}{n} \mathbf{J}_n$.

- (a) Halle la distribución de Q_k , $k = 1, 2$.
 (b) ¿Son Q_1 y Q_2 independientes?

Inferencia en el Modelo Lineal

En este capítulo se describe la inferencia en modelos lineales. Primeramente introducimos algunas definiciones y supuestos en los que se basan los modelos de regresión.

2.1. Definición de un modelo lineal

DEFINICIÓN 2.1. Considere la variable aleatoria Y , decimos que sigue un *modelo lineal* si

$$E(Y) = \sum_{j=1}^p x_j \beta_j = \mathbf{x}^\top \boldsymbol{\beta},$$

donde $\mathbf{x} = (x_1, \dots, x_p)^\top$ representa un vector de p *variables regresoras*, mientras que $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ denota un vector de parámetros desconocidos, conocidos como *coeficientes de regresión*.

Suponga que tenemos n observaciones recolectadas desde Y , entonces podemos considerar el modelo

$$E(Y_i) = \mathbf{x}_i^\top \boldsymbol{\beta}, \quad i = 1, \dots, n.$$

Es conveniente escribir lo anterior como:

$$E(\mathbf{Y}) = \begin{pmatrix} \mathbf{x}_1^\top \boldsymbol{\beta} \\ \vdots \\ \mathbf{x}_n^\top \boldsymbol{\beta} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} \boldsymbol{\beta} = \mathbf{X} \boldsymbol{\beta},$$

donde $\mathbf{X} = (x_{ij}) \in \mathbb{R}^{n \times p}$ se denomina *matriz de diseño*.

OBSERVACIÓN 2.2. Cuando

$$\text{Cov}(\mathbf{Y}) = \sum_{r=1}^K \phi_r \boldsymbol{\Sigma}_r = \boldsymbol{\Sigma}(\boldsymbol{\phi}),$$

donde $\boldsymbol{\phi} = (\phi_1, \dots, \phi_K)^\top$ son parámetros desconocidos y $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K$ son matrices (simétricas) conocidas, decimos que \mathbf{Y} sigue un *modelo lineal general*.

Típicamente, asumiremos que Y_1, \dots, Y_n son independientes con varianza constante, en cuyo caso

$$\text{Cov}(\mathbf{Y}) = \sigma^2 \mathbf{I},$$

y decimos que \mathbf{Y} sigue un *modelo lineal simple*.

DEFINICIÓN 2.3. Se dice que el vector \mathbf{Y} sigue un *modelo lineal* si

$$E(\mathbf{Y}) = \mathbf{X} \boldsymbol{\beta}, \quad \text{Cov}(\mathbf{Y}) = \sum_{r=1}^K \phi_r \boldsymbol{\Sigma}_r.$$

Note que la definición anterior puede ser expresada de forma equivalente como:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

con

$$\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}, \quad \text{Cov}(\boldsymbol{\epsilon}) = \sum_{r=1}^K \phi_r \boldsymbol{\Sigma}_r.$$

OBSERVACIÓN 2.4. Los supuestos de momentos dados en la Definición 2.3 suelen ser llamados *condiciones de Gauss-Markov*. Aunque es usual que sean expresados en términos del modelo lineal simple, como

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

con

$$\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}, \quad \text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n.$$

DEFINICIÓN 2.5. Se dice que el vector \mathbf{Y} sigue un *modelo lineal normal* si

$$\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\phi)),$$

o bien

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathbf{N}_n(\mathbf{0}, \boldsymbol{\Sigma}(\phi)).$$

Mientras que sigue un *modelo normal simple*, si

$$\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}). \quad (2.1)$$

OBSERVACIÓN 2.6. Se debe destacar que el modelo en (2.1) puede ser escrito como

$$Y_i \stackrel{\text{ind}}{\sim} \mathbf{N}_1(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2), \quad i = 1, \dots, n. \quad (2.2)$$

En efecto, la inferencia estadística para los modelos definidos en Ecuaciones (2.1) y (2.2) son equivalentes.¹

SUPUESTO 1. *El modelo lineal descrito por la ecuación:*

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

tiene los siguientes supuestos:

- A1: \mathbf{X} es una matriz no aleatoria $n \times p$ con $n > p$.
- A2: La matriz \mathbf{X} tiene rango p , es decir, \mathbf{X} es rango columna completo.
- A3: El vector aleatorio n -dimensional \mathbf{Y} tiene elementos que son observables.
- A4: El vector aleatorio no observable $\boldsymbol{\epsilon}$ satisface

$$\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}, \quad \text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}, \quad \sigma^2 > 0.$$

El Supuesto A4 puede ser re-establecido para incorporar la suposición de normalidad, esto es,

SUPUESTO 2. *Considere:*

A4*: *El vector aleatorio $\boldsymbol{\epsilon}$ satisface $\boldsymbol{\epsilon} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, o equivalentemente*

$$\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n).$$

¹Es decir, **bajo normalidad** estos dos modelos son equivalentes. En efecto, esto **no** es verdad en general.

2.2. Estimación de parámetros en el modelo de regresión lineal

Primeramente abordaremos la estimación máximo verosímil bajo normalidad, revisaremos propiedades de los estimadores y abordaremos la conexión con el método de mínimos cuadrados. En efecto, es fácil notar que la función de verosimilitud proveniente del modelo $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$, adopta la forma:

$$L(\boldsymbol{\theta}) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2\right),$$

con $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$. De este modo, la función de log-verosimilitud es dada por

$$\begin{aligned}\ell(\boldsymbol{\theta}) &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2}\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\ &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2}Q(\boldsymbol{\beta}),\end{aligned}$$

donde

$$Q(\boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2,$$

se denomina *suma de cuadrados del error*. Diferenciando con relación a $\boldsymbol{\beta}$ y σ^2 , obtenemos

$$\mathbf{d}_\beta \ell(\boldsymbol{\theta}) = -\frac{1}{2\sigma^2} \mathbf{d}_\beta \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \frac{1}{\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{X} \mathbf{d}\boldsymbol{\beta},$$

y

$$\mathbf{d}_{\sigma^2} \ell(\boldsymbol{\theta}) = -\frac{n}{2\sigma^2} \mathbf{d}\sigma^2 + \frac{1}{2\sigma^4} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 \mathbf{d}\sigma^2.$$

Es decir,²

$$\begin{aligned}\frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} &= \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}), \\ \frac{\partial \ell(\boldsymbol{\theta})}{\partial \sigma^2} &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2.\end{aligned}$$

Desde la condición de primer orden, tenemos las *ecuaciones de verosimilitud*:

$$\begin{aligned}\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) &= \mathbf{0} \\ n\hat{\sigma}^2 - \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 &= 0.\end{aligned}$$

Resolviendo con relación a $\boldsymbol{\beta}$ obtenemos las *ecuaciones normales*

$$\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{Y},$$

dado que $\text{rg}(\mathbf{X}) = \text{rg}(\mathbf{X}^\top \mathbf{X}) = p$, entonces el sistema anterior admite solución única dada por:³

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}, \\ \hat{\sigma}^2 &= \frac{1}{n} Q(\hat{\boldsymbol{\beta}}) = \frac{1}{n} \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2.\end{aligned}$$

Se define el vector de *valores predichos* como

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \mathbf{H}\mathbf{Y},$$

donde $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

²Usando el primer Teorema de identificación dado en Magnus y Neudecker (2007), pag. 98.

³Note que $\hat{\boldsymbol{\beta}} = \mathbf{X}^+ \mathbf{Y}$ con $\mathbf{X}^+ = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ la inversa Moore-Penrose de \mathbf{X} .

Es fácil notar que \mathbf{H} es simétrica e idempotente, en cuyo caso

$$\text{rg}(\mathbf{H}) = \text{tr}(\mathbf{H}) = \text{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \text{tr}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}) = p.$$

El vector de diferencias entre \mathbf{Y} y $\hat{\mathbf{Y}}$ se denomina el *vector de residuos*, es decir

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}.$$

Además, tenemos que $\mathbf{I} - \mathbf{H}$ también es simétrica e idempotente. Con esta notación podemos escribir

$$\begin{aligned} Q(\hat{\beta}) &= \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 = \mathbf{e}^\top \mathbf{e} = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H})\mathbf{Y} \\ &= \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \hat{\beta})^2, \end{aligned}$$

que es conocido como *suma de cuadrados residual*.

RESULTADO 2.7. *Considere el modelo:*

$$\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n),$$

con los supuestos A1 a A4*. Entonces, tenemos que:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \sim \mathbf{N}_p(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}), \quad (2.3)$$

de donde sigue que $\hat{\beta}$ es insesgado y $\text{Cov}(\hat{\beta}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$. Además, la variable aleatoria

$$\frac{Q(\hat{\beta})}{\sigma^2} = \frac{\mathbf{Y}^\top (\mathbf{I} - \mathbf{H})\mathbf{Y}}{\sigma^2} \sim \chi^2(n - p).$$

Finalmente $\hat{\beta}$ y $s^2 = Q(\hat{\beta})/(n - p)$ son independientes.

DEMOSTRACIÓN. Notando que $\hat{\beta}$ es una transformación lineal de un vector aleatorio normal y

$$\begin{aligned} \mathbb{E}(\hat{\beta}) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{Y}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\beta = \beta, \\ \text{Cov}(\hat{\beta}) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Cov}(\mathbf{Y}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}, \end{aligned}$$

y el resultado en Ecuación (2.3) sigue. Por otro lado,

$$\frac{Q(\hat{\beta})}{\sigma^2} = \frac{(n - p)s^2}{\sigma^2} = \frac{\mathbf{Y}^\top (\mathbf{I} - \mathbf{H})\mathbf{Y}}{\sigma^2},$$

sigue una distribución chi-cuadrado pues $\mathbf{I} - \mathbf{H}$ es matriz idempotente con

$$\text{rg}(\mathbf{I} - \mathbf{H}) = \text{tr}(\mathbf{I} - \mathbf{H}) = \text{tr}(\mathbf{I}) - \text{tr}(\mathbf{H}) = n - p.$$

Además, el parámetro de no centralidad es dado por

$$\lambda = \frac{1}{2\sigma^2} \beta^\top \mathbf{X}^\top (\mathbf{I} - \mathbf{H}) \mathbf{X} \beta = 0.$$

En efecto,

$$\mathbf{H}\mathbf{X} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} = \mathbf{X} \quad \Rightarrow \quad (\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}.$$

La independencia entre $\hat{\beta}$ y s^2 sigue desde $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$, lo que concluye la prueba. \square

OBSERVACIÓN 2.8. Es fácil notar que

$$\begin{aligned} \mathbb{E}\{Q(\hat{\beta})\} &= \mathbb{E}\{\mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}\} = \sigma^2 \text{tr}(\mathbf{I} - \mathbf{H}) + \beta^\top \mathbf{X}^\top (\mathbf{I} - \mathbf{H}) \mathbf{X} \beta \\ &= \sigma^2(n - p), \end{aligned}$$

es decir,

$$\mathbb{E}(\hat{\sigma}^2) = \left(\frac{n-p}{n}\right) \sigma^2,$$

lo que permite sugerir el *estimador insesgado*:

$$s^2 = \frac{1}{n-p} \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2.$$

RESULTADO 2.9. *El vector de valores predichos y el vector de residuos son independientemente distribuidos como*

$$\hat{\mathbf{Y}} \sim \mathcal{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{H}), \quad \mathbf{e} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H})).$$

DEMOSTRACIÓN. Considere

$$\begin{pmatrix} \hat{\mathbf{Y}} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{H} \\ \mathbf{I} - \mathbf{H} \end{pmatrix} \mathbf{Y},$$

luego,

$$\mathbb{E} \begin{pmatrix} \hat{\mathbf{Y}} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{H} \\ \mathbf{I} - \mathbf{H} \end{pmatrix} \mathbf{X}\beta = \begin{pmatrix} \mathbf{H}\mathbf{X}\beta \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{X}\beta \\ \mathbf{0} \end{pmatrix}.$$

Además,

$$\begin{aligned} \text{Cov} \begin{pmatrix} \hat{\mathbf{Y}} \\ \mathbf{e} \end{pmatrix} &= \sigma^2 \begin{pmatrix} \mathbf{H} \\ \mathbf{I} - \mathbf{H} \end{pmatrix} (\mathbf{H}^\top, (\mathbf{I} - \mathbf{H})^\top) \\ &= \sigma^2 \begin{pmatrix} \mathbf{H}\mathbf{H}^\top & \mathbf{H}(\mathbf{I} - \mathbf{H})^\top \\ (\mathbf{I} - \mathbf{H})\mathbf{H}^\top & (\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H})^\top \end{pmatrix} \\ &= \sigma^2 \begin{pmatrix} \mathbf{H}^2 & \mathbf{0} \\ \mathbf{0} & (\mathbf{I} - \mathbf{H})^2 \end{pmatrix} \\ &= \sigma^2 \begin{pmatrix} \mathbf{H} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \mathbf{H} \end{pmatrix}, \end{aligned}$$

lo que permite establecer el resultado. \square

OBSERVACIÓN 2.10. Debemos resaltar que \mathbf{H} y $\mathbf{I} - \mathbf{H}$ son matrices de rango incompleto y por tanto $\hat{\mathbf{Y}}$ y \mathbf{e} siguen distribuciones normales *singulares*.

A continuación vamos a suponer que el vector de respuestas y la matriz de diseño dependen del tamaño muestral, n . Considere el siguiente problema

$$\min_{\beta} Q_n(\beta) = \min_{\beta} \|\mathbf{Y}_n - \mathbf{X}_n \beta\|^2,$$

cuya solución, $\hat{\beta}_n$ es llamado estimador *mínimos cuadrados (LS)*, dado por:

$$\hat{\beta}_n = (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{Y}_n,$$

Es habitual asumir el modelo

$$\mathbf{Y}_n = \mathbf{X}_n \beta + \epsilon_n,$$

con $E(\epsilon_n) = \mathbf{0}$ y $\text{Cov}(\epsilon_n) = \sigma^2 \mathbf{I}_n$. Esto lleva a

$$\begin{aligned} E(\hat{\beta}_n) &= (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top E(\mathbf{X}_n \beta + \epsilon_n) = \beta \\ \text{Cov}(\hat{\beta}_n) &= \text{Cov}((\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top (\mathbf{X}_n \beta + \epsilon_n)) = \sigma^2 (\mathbf{X}_n^\top \mathbf{X}_n)^{-1}. \end{aligned}$$

SUPUESTO 3. *Considere el modelo,*

$$\mathbf{Y}_n = \mathbf{X}_n \beta + \epsilon_n, \quad (2.4)$$

y suponga las condiciones:

B1: $E(\epsilon_n) = \mathbf{0}$, $\text{Cov}(\epsilon_n) = \sigma^2 \mathbf{I}$.

B2: Sea $h_{kk} = \mathbf{x}_{k,n}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{k,n}$ el k -ésimo elemento de la diagonal de \mathbf{H}_n y considere

$$\max_{1 \leq k \leq n} h_{kk} \rightarrow 0, \quad \text{conforme } n \rightarrow \infty.$$

B3: $\lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}_n^\top \mathbf{X}_n = \mathbf{K}$ es una matriz no singular (no estocástica).

Usando el supuesto B3, tenemos

$$\lim_{n \rightarrow \infty} \text{Cov}(\hat{\beta}_n) = \sigma^2 \lim_{n \rightarrow \infty} \frac{1}{n} \left(\frac{\mathbf{X}_n^\top \mathbf{X}_n}{n} \right)^{-1} = \sigma^2 \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{K}^{-1} = \mathbf{0}.$$

Esto implica que $\hat{\beta}_n \xrightarrow{2\text{nd}} \beta$. Es decir, $\hat{\beta}_n$ es un estimador consistente de β .

RESULTADO 2.11 (Distribución asintótica del estimador LS). *Considere el modelo (2.4) bajo los supuestos B1 a B3. Entonces*

$$\sqrt{n}(\hat{\beta}_n - \beta) \xrightarrow{D} N_p(\mathbf{0}, \sigma^2 \mathbf{K}^{-1}).$$

DEMOSTRACIÓN. Ver Sen y Singer (1993), pág. 279. \square

RESULTADO 2.12 (Teorema de Gauss-Markov). *Suponga $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ y sea $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ el estimador mínimos cuadrados. Asuma el supuesto B1. El estimador $\hat{\gamma} = \mathbf{h}^\top \hat{\beta}$ de $\gamma = \mathbf{h}^\top \beta$ es el mejor estimador lineal e insesgado (BLUE).*

DEMOSTRACIÓN. Sea $\mathbf{c}^\top \mathbf{Y}$ cualquier otro estimador lineal e insesgado de $\gamma = \mathbf{h}^\top \beta$. Dado que $\mathbf{c}^\top \mathbf{Y}$ es insesgado, sigue que

$$\mathbf{h}^\top \beta = E(\mathbf{c}^\top \mathbf{Y}) = \mathbf{c}^\top E(\mathbf{Y}) = \mathbf{c}^\top \mathbf{X} \beta, \quad \forall \beta,$$

luego, tenemos

$$\mathbf{c}^\top \mathbf{X} = \mathbf{h}^\top.$$

Ahora,

$$\text{var}(\mathbf{c}^\top \mathbf{Y}) = \mathbf{c}^\top \text{Cov}(\mathbf{Y}) \mathbf{c} = \sigma^2 \mathbf{c}^\top \mathbf{c}, \quad (2.5)$$

mientras que

$$\text{var}(\mathbf{h}^\top \hat{\beta}) = \mathbf{h}^\top \text{Cov}(\hat{\beta}) \mathbf{h} = \sigma^2 \mathbf{h}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{h} = \sigma^2 \mathbf{c}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{c}, \quad (2.6)$$

desde (2.5) y (2.6), tenemos

$$\begin{aligned} \text{var}(\mathbf{c}^\top \mathbf{Y}) - \text{var}(\mathbf{h}^\top \hat{\beta}) &= \sigma^2 (\mathbf{c}^\top \mathbf{c} - \mathbf{c}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{c}) \\ &= \sigma^2 \mathbf{c}^\top (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}) \mathbf{c} \geq 0, \end{aligned}$$

el resultado sigue pues $\mathbf{I} - \mathbf{H}$ es semidefinida positiva. \square

2.3. Aspectos numéricos de estimación LS en regresión lineal

El estimador mínimos cuadrados (o de máximo verosimilitud bajo normalidad) para el modelo en lineal con los Supuestos A1-A4, puede ser expresado como la solución del problema:

$$\min_{\beta \in \mathbb{R}^p} Q(\beta), \quad \text{con} \quad Q(\beta) = \|Y - X\beta\|_2^2,$$

lo que lleva a las ecuaciones de estimación $X^\top(Y - X\hat{\beta}) = \mathbf{0}$, o bien

$$X^\top X \hat{\beta} = X^\top Y. \quad (2.7)$$

Métodos habituales para obtener $\hat{\beta}$, son:

- Métodos directos basados en la factorización Cholesky, el operador Sweep (Goodnight, 1979) y descomposiciones QR y SVD.
- Un poco menos común en regresión, es el uso del método gradientes conjugados (CG) (McIntosh, 1982).

OBSERVACIÓN 2.13. Algunas características de los procedimientos descritos anteriormente son relevantes de ser destacadas, a saber:

- Cholesky y Sweep requieren formar las matrices:

$$X^\top X, X^\top Y, \quad \text{y} \quad \begin{pmatrix} X^\top X & X^\top Y \\ Y^\top X & Y^\top Y \end{pmatrix},$$

respectivamente.

- QR y SVD descomponen la matriz de diseño X y resuelven sistemas lineales (triangular y diagonal, respectivamente) mucho más pequeños ($n \gg p$).
- Una implementación cuidadosa de CG sólo requiere productos matriz-vector/operaciones entre vectores y tan sólo $4p$ posiciones para almacenamiento.
- Existe código confiable y con excelente desempeño para álgebra lineal numérica en las bibliotecas BLAS, LAPACK, rutinas que pueden ser invocadas desde (por ejemplo) R y Matlab.

A continuación introducimos algunas ideas sobre los procedimientos numéricos que serán aplicados al modelo de regresión lineal. Primeramente, considere el siguiente resultado

RESULTADO 2.14 (Factorización Cholesky). *Sea $A \in \mathbb{R}^{p \times p}$ es matriz simétrica y definida positiva, entonces existe una única matriz triangular superior $G \in \mathbb{R}^{p \times p}$ con elementos diagonales positivos tal que*

$$A = G^\top G$$

Note que si usamos la factorización Cholesky para resolver el sistema $Ax = b$. Entonces debemos resolver los sistemas triangulares

$$G^\top z = b, \quad \text{y} \quad Gx = z.$$

En efecto,

$$Ax = (G^\top G)x = G^\top(Gx) = G^\top z = b.$$

Algoritmo 1: Factorización Cholesky**Entrada:** Matriz $\mathbf{A} \in \mathbb{R}^{p \times p}$.**Salida :** Factor Cholesky $\mathbf{G} \in \mathbb{R}^{p \times p}$.

```

1 begin
2    $g_{11} = \sqrt{a_{11}}$ .
3   for  $j = 2$  to  $p$  do
4      $g_{1j} = a_{1j}/g_{11}$ .
5   end
6   for  $i = 2$  to  $p$  do
7      $g_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} g_{ki}^2}$ ,
8     for  $j = i + 1$  to  $n$  do
9        $g_{ij} = (a_{ij} - \sum_{k=1}^{i-1} g_{ki}g_{kj})/g_{ii}$ 
10    end
11  end
12 end

```

Sea $RSS = Q(\hat{\beta})$ la suma de cuadrados residuales y notando que

$$RSS = \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 = \mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{X}\hat{\beta} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y} + \hat{\beta}^\top \mathbf{X}^\top \mathbf{X}\hat{\beta}.$$

Tenemos,

$$\mathbf{Y}^\top \mathbf{X}\hat{\beta} = \mathbf{Y}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \mathbf{Y}^\top \mathbf{H}^2 \mathbf{Y} = \hat{\beta}^\top \hat{\mathbf{Y}} = \hat{\beta}^\top \mathbf{X}^\top \mathbf{X}\hat{\beta}.$$

Es decir, podemos escribir:

$$RSS = \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 = \mathbf{Y}^\top \mathbf{Y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{X}\hat{\beta}.$$

Podemos resolver (2.7) usando la descomposición Cholesky, de

$$\mathbf{X}^\top \mathbf{X} = \mathbf{U}^\top \mathbf{U},$$

con \mathbf{U} matrix triangular superior. De este modo, debemos resolver los sistemas triangulares:

$$\mathbf{U}^\top \mathbf{z} = \mathbf{X}^\top \mathbf{Y}, \quad \text{y} \quad \mathbf{U}\hat{\beta} = \mathbf{z}.$$

Mientras que para obtener s^2 , consideramos:

$$RSS = \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 = \mathbf{Y}^\top \mathbf{Y} - \mathbf{z}^\top \mathbf{z}.$$

Adicionalmente,

$$\mathbf{U}^{-1} \mathbf{U}^{-\top} = (\mathbf{X}^\top \mathbf{X})^{-1},$$

es proporcional a la matriz de covarianza de $\hat{\beta}$.

OBSERVACIÓN 2.15. Invertiendo \mathbf{U} in-place, esto es, haciendo

$$\mathbf{U} \leftarrow \mathbf{U}^{-1},$$

sigue que $(\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{U}\mathbf{U}^\top$ lo que permite ahorrar espacio de almacenamiento y puede ser eficientemente calculado usando rutinas desde BLAS.

DEFINICIÓN 2.16 (Operador Sweep). Sea $\mathbf{A} = (a_{ij})$ matriz cuadrada $p \times p$, aplicando el operador *Sweep* sobre el k -ésimo elemento diagonal de \mathbf{A} ($a_{kk} \neq 0$) permite obtener la matriz \mathbf{B} , definida como:

$$\begin{aligned} b_{kk} &= \frac{1}{a_{kk}}, \\ b_{ik} &= -\frac{a_{ik}}{a_{kk}}, & i \neq k, \\ b_{kj} &= \frac{a_{kj}}{a_{kk}}, & j \neq k, \\ b_{ij} &= a_{ij} - \frac{a_{ik}a_{kj}}{a_{kk}}, & i, j \neq k, \end{aligned}$$

y escribimos $\mathbf{B} = \text{Sweep}(k)\mathbf{A}$.

PROPIEDAD 2.17. El operador Sweep disfruta de las siguientes propiedades:

- (i) $\text{Sweep}(k)\text{Sweep}(k)\mathbf{A} = \mathbf{A}$.
- (ii) $\text{Sweep}(k)\text{Sweep}(r)\mathbf{A} = \text{Sweep}(r)\text{Sweep}(k)\mathbf{A}$.
- (iii) $\mathbf{A}^{-1} = \prod_{i=1}^n \text{Sweep}(i)\mathbf{A}$.

Debemos destacar que, si \mathbf{A} es matriz simétrica, el operador Sweep *preserva la simetría* de \mathbf{A} . Existen varias definiciones ligeramente diferentes del operador Sweep y más importante, es conocido que problemas de *inestabilidad* pueden ocurrir cuando algún a_{kk} es cercano a cero.

Considere $\mathbf{A} \in \mathbb{R}^{p \times p}$ matriz particionada como:

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix},$$

donde $\mathbf{A}_{11} \in \mathbb{R}^{r \times r}$ ($r < p$). Suponga que se aplica el operador Sweep sobre los elementos diagonales de \mathbf{A}_{11} . De este modo,

$$\mathbf{B} = \prod_{i=1}^r \text{Sweep}(i)\mathbf{A} = \begin{pmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{pmatrix},$$

con

$$\begin{aligned} \mathbf{B}_{11} &= \mathbf{A}_{11}^{-1}, & \mathbf{B}_{12} &= \mathbf{A}_{11}^{-1}\mathbf{A}_{12}, \\ \mathbf{B}_{21} &= -\mathbf{A}_{21}\mathbf{A}_{11}^{-1}, & \mathbf{B}_{22} &= \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}. \end{aligned}$$

Este resultado nos permite utilizar el operador Sweep en problemas de regresión. Considere

$$\mathbf{Z} = (\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^{n \times (p+1)},$$

luego

$$\mathbf{Z}^\top \mathbf{Z} = \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{Y} \\ \mathbf{Y}^\top \mathbf{X} & \mathbf{Y}^\top \mathbf{Y} \end{pmatrix} \in \mathbb{R}^{(p+1) \times (p+1)}.$$

que corresponde a una matriz cuadrada de orden $(p+1) \times (p+1)$.

Aplicando el operador Sweep sobre los primeros p elementos diagonales de $\mathbf{Z}^\top \mathbf{Z}$, obtenemos:

$$\begin{aligned} \mathbf{B} &= \prod_{i=1}^p \text{Sweep}(i) \mathbf{Z}^\top \mathbf{Z} \\ &= \begin{pmatrix} (\mathbf{X}^\top \mathbf{X})^{-1} & (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \\ -\mathbf{Y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} & \mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \end{pmatrix} \\ &= \begin{pmatrix} (\mathbf{X}^\top \mathbf{X})^{-1} & \hat{\boldsymbol{\beta}} \\ -\hat{\boldsymbol{\beta}}^\top & \text{RSS} \end{pmatrix}. \end{aligned}$$

es decir, este procedimiento nos permite calcular todos los elementos necesarios para llevar a cabo la estimación en un modelo de regresión lineal.

En el contexto de regresión, es usual tener que $n \gg p$ y uno de los procedimientos preferidos para llevar a cabo la estimación mínimos cuadrados esta basado en la descomposición ortogonal-triangular (QR) de la matriz de diseño \mathbf{X} . Considere la siguiente definición,

DEFINICIÓN 2.18 (Descomposición QR). Sea $\mathbf{A} \in \mathbb{R}^{n \times p}$, entonces existe $\mathbf{Q} \in \mathcal{O}_n$ y $\mathbf{R} \in \mathbb{R}^{n \times p}$, tal que

$$\mathbf{A} = \mathbf{Q}\mathbf{R}, \quad \mathbf{R} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix}$$

donde $\mathbf{R}_1 \in \mathbb{R}^{p \times p}$ matriz triangular superior, aquí suponemos que $n \geq p$.

En efecto, es fácil notar que si $\mathbf{A} = \mathbf{Q}\mathbf{R}$, entonces

$$\mathbf{A}^\top \mathbf{A} = \mathbf{R}^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{R} = \mathbf{R}^\top \mathbf{R} = \mathbf{R}_1^\top \mathbf{R}_1,$$

y \mathbf{R}_1 corresponde al factor Cholesky de $\mathbf{A}^\top \mathbf{A}$. Este aspecto es relevante pues no es necesario formar la matriz de productos cruzados $\mathbf{A}^\top \mathbf{A}$ para obtener el factor Cholesky \mathbf{R}_1 . Antes de describir brevemente el algoritmo para obtener la descomposición QR recordamos algunas propiedades fundamentales de las matrices ortogonales:

- $\mathbf{Q}\mathbf{Q}^\top = \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$.
- $\langle \mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{y} \rangle = \mathbf{x}^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{y} = \mathbf{x}^\top \mathbf{y} = \langle \mathbf{x}, \mathbf{y} \rangle$.
- $\|\mathbf{Q}\mathbf{x}\| = \|\mathbf{x}\|$.
- Si $\mathbf{B} = \mathbf{Q}^\top \mathbf{A} \mathbf{Q}$, entonces \mathbf{A} y \mathbf{B} tienen los mismos valores propios para \mathbf{Q} matriz ortogonal.

Existen diversas variantes del algoritmo para implementar la descomposición QR. A continuación veremos una basada en *transformaciones Householder*. Primeramente, considere el siguiente problema

PROBLEMA 2.19. Para $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{x} \neq \mathbf{0}$, hallar una matriz ortogonal $\mathbf{M} \in \mathbb{R}^{n \times n}$ tal que

$$\mathbf{M}^\top \mathbf{x} = \|\mathbf{x}\| \mathbf{e}_1,$$

donde $\mathbf{e}_1 = (1, 0, \dots, 0)^\top$ denota el primer vector unidad.

DEFINICIÓN 2.20 (Reflexión). Sea \mathbf{u} y \mathbf{v} vectores ortonormales y \mathbf{x} vector generado por \mathbf{u} y \mathbf{v} . Entonces

$$\mathbf{x} = c_1\mathbf{u} + c_2\mathbf{v},$$

para escalares c_1, c_2 . El vector

$$\tilde{\mathbf{x}} = -c_1\mathbf{u} + c_2\mathbf{v},$$

el llamado una *reflexión* de \mathbf{x} a través de la línea definida por el vector \mathbf{v} (o \mathbf{u}^\perp).

DEFINICIÓN 2.21 (Transformación Householder). Sea $\mathbf{x} = c_1\mathbf{u} + c_2\mathbf{v}$, con \mathbf{u} y \mathbf{v} vectores generadores de \mathbf{x} y considere la matriz

$$\mathbf{H} = \mathbf{I} - \lambda\mathbf{u}\mathbf{u}^\top, \quad \lambda = 2/\mathbf{u}^\top\mathbf{u}.$$

Note que $\mathbf{H}\mathbf{x} = \tilde{\mathbf{x}}$, es decir \mathbf{H} es un reflector.

El objetivo es determinar una matriz \mathbf{M} basado en reflexiones Householder. Debemos destacar que la *transformación Householder* satisface las siguientes propiedades:

- $\mathbf{H}\mathbf{u} = -\mathbf{u}$.
- $\mathbf{H}\mathbf{v} = \mathbf{v}$ para cualquier \mathbf{v} ortogonal a \mathbf{u} .
- $\mathbf{H}^\top = \mathbf{H}$.
- $\mathbf{H}^{-1} = \mathbf{H}^\top$.

OBSERVACIÓN 2.22. La operación $\mathbf{H}\mathbf{x}$ puede ser obtenida usando una operación axpy.⁴ En efecto,

$$\mathbf{H}\mathbf{x} = (\mathbf{I} - \lambda\mathbf{u}\mathbf{u}^\top)\mathbf{x} = \mathbf{x} - \alpha\mathbf{u}, \quad \alpha = \lambda\mathbf{u}^\top\mathbf{x}.$$

Algoritmo 2: Descomposición QR

Entrada: Matriz $\mathbf{A} \in \mathbb{R}^{n \times p}$.

Salida : Factores \mathbf{Q} y \mathbf{R} , matrices ortogonal y triangular superior, respectivamente.

```

1 begin
2   Hacer  $\mathbf{Q} = \mathbf{I}_n$  y  $\mathbf{R} = \mathbf{A}$ 
3   for  $i = 1$  to  $p$  do
4      $\mathbf{x} = (R_{1i}, \dots, R_{pi})^\top$ 
5      $\mathbf{Q}_i = \begin{pmatrix} \mathbf{I}_{i-1} & \mathbf{0} \\ \mathbf{0} & M(\mathbf{x}) \end{pmatrix}$ 
6     /*  $M(\mathbf{x})$  obtenido usando reflexiones Householder */
7      $\mathbf{Q} = \mathbf{Q}_i\mathbf{Q}$ 
8      $\mathbf{R} = \mathbf{Q}_i\mathbf{R}$ 
9   end
10   $\mathbf{Q} = \mathbf{Q}^\top$ 
11   $\mathbf{R} = (R_{ij})$  para  $i, j = 1, \dots, p$ .
12 end
```

⁴Corresponde a una actualización del tipo: $\mathbf{y} \leftarrow \alpha\mathbf{x} + \mathbf{y}$.

La descomposición QR de una matriz $\mathbf{A} \in \mathbb{R}^{n \times p}$ ($n > p$), puede ser construída a través de una secuencia de matrices $\mathbf{Q}_1, \dots, \mathbf{Q}_p$ tales que

$$\mathbf{Q}_p \cdots \mathbf{Q}_1 \mathbf{A} = \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix},$$

donde $\mathbf{Q}_1, \dots, \mathbf{Q}_p$ son *todas* ortogonales. De este modo,

$$\mathbf{A} = \mathbf{Q}_1^\top \cdots \mathbf{Q}_p^\top \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} = \mathbf{Q} \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}.$$

El Algoritmo 2 permite obtener la descomposición QR de $\mathbf{A} \in \mathbb{R}^{n \times p}$ usando transformaciones Householder.⁵ En este contexto, $\mathbf{M}(\mathbf{x})$ corresponde a la matriz ortogonal desde el Problema 2.19 basada en un vector \mathbf{x} . Debemos destacar que una implementación eficiente del Algoritmo 2 solamente requiere almacenar la información mínima para formar las matrices $\mathbf{Q}_1, \dots, \mathbf{Q}_p$ y la propia matriz triangular \mathbf{R}_1 en la propia matriz \mathbf{A} , lo que permite un ahorro desde el punto de vista de almacenamiento.

Para el problema de regresión, considere la descomposición QR de \mathbf{X} , como:

$$\mathbf{X} = \mathbf{Q}\mathbf{R}, \quad \mathbf{R} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix},$$

con $\mathbf{R}_1 \in \mathbb{R}^{p \times p}$ matriz triangular superior ($n > p$). Si $\text{rg}(\mathbf{X}) = p$, entonces \mathbf{R}_1 es no singular. Además, considere la transformación:

$$\mathbf{Q}^\top \mathbf{Y} = \mathbf{c}, \quad \mathbf{c} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix}.$$

La descomposición QR de \mathbf{X} permite re-escribir la función objetivo asociada al problema mínimos cuadrados, como:

$$\begin{aligned} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 &= \|\mathbf{Q}^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\|^2 = \|\mathbf{Q}^\top \mathbf{Y} - \mathbf{Q}^\top \mathbf{Q}\mathbf{R}\boldsymbol{\beta}\|^2 \\ &= \|\mathbf{c} - \mathbf{R}\boldsymbol{\beta}\|^2, \end{aligned}$$

Es fácil notar que:

$$\|\mathbf{c} - \mathbf{R}\boldsymbol{\beta}\|^2 = \|\mathbf{c}_1 - \mathbf{R}_1\boldsymbol{\beta}\|^2 + \|\mathbf{c}_2\|^2.$$

Esto lleva a escribir el estimador de mínimos cuadrados $\hat{\boldsymbol{\beta}}$ como solución del sistema triangular:

$$\mathbf{R}_1 \hat{\boldsymbol{\beta}} = \mathbf{c}_1.$$

El mínimo de la función objetivo está dado por $\|\mathbf{c}_2\|^2$, lo que permite calcular el estimador insesgado de σ^2 como

$$s^2 = \frac{1}{n-p} \|\mathbf{c}_2\|^2 = \frac{\text{RSS}}{n-p},$$

Finalmente,

$$\mathbf{X}^\top \mathbf{X} = (\mathbf{R}_1^\top, \mathbf{0}) \mathbf{Q}^\top \mathbf{Q} \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix} = \mathbf{R}_1^\top \mathbf{R}_1,$$

lo que ofrece un procedimiento eficiente para obtener la matriz de covarianza de $\hat{\boldsymbol{\beta}}$, dado por

$$\text{Cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{R}_1^\top \mathbf{R}_1)^{-1} = \sigma^2 \mathbf{R}_1^{-1} \mathbf{R}_1^{-\top}.$$

⁵Otro método popular para obtener la descomposición QR es usando rotaciones Givens.

DEFINICIÓN 2.23 (Descomposición Valor Singular). Sea $\mathbf{A} \in \mathbb{R}^{n \times p}$ con $\text{rg}(\mathbf{A}) = r$, entonces existen matrices $\mathbf{U} \in \mathcal{O}_n$, $\mathbf{V} \in \mathcal{O}_p$, tal que

$$\mathbf{A} = \mathbf{U} \begin{pmatrix} \mathbf{D}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V}^\top,$$

donde $\mathbf{D}_r = \text{diag}(\delta_1, \dots, \delta_r)$ con $\delta_1 \geq \delta_2 \geq \dots \geq \delta_r > 0$, que son llamados valores singulares de \mathbf{A} .

La Descomposición Valor Singular (SVD) para $\mathbf{A} \in \mathbb{R}^{n \times p}$ con $\text{rg}(\mathbf{A}) = r$ puede ser escrita como:

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^\top,$$

con $\mathbf{U} \in \mathbb{R}^{n \times p}$ tal que $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_p$, $\mathbf{D} = \text{diag}(\delta_1, \dots, \delta_r)$ y $\mathbf{V} \in \mathcal{O}_p$.

Para el contexto de regresión lineal, considere la SVD de \mathbf{X} ,

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top,$$

donde $\mathbf{U} \in \mathbb{R}^{n \times p}$ tal que $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_p$, $\mathbf{D} = \text{diag}(\delta_1, \dots, \delta_p)$ y $\mathbf{V} \in \mathcal{O}_p$. De este modo, podemos escribir el modelo de regresión lineal como:

$$\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{U} \mathbf{D} \boldsymbol{\alpha} + \boldsymbol{\epsilon},$$

con $\boldsymbol{\alpha} = \mathbf{V}^\top \boldsymbol{\beta}$. Haciendo $\mathbf{Z} = \mathbf{U}^\top \mathbf{Y}$, tenemos el modelo en *forma canónica*:

$$\mathbf{Z} = \mathbf{D} \boldsymbol{\alpha} + \boldsymbol{\eta}, \quad \boldsymbol{\eta} = \mathbf{U}^\top \boldsymbol{\epsilon},$$

donde

$$\mathbf{E}(\boldsymbol{\eta}) = \mathbf{0}, \quad \text{Cov}(\boldsymbol{\eta}) = \sigma^2 \mathbf{U}^\top \mathbf{U} = \sigma^2 \mathbf{I}_p.$$

Es decir, podemos el modelo canónico satisface las condiciones A1 a A4. Esto lleva al estimador LS de $\boldsymbol{\alpha}$ en el modelo canónico,

$$\hat{\boldsymbol{\alpha}} = \mathbf{D}^{-1} \mathbf{Z}, \quad \Rightarrow \quad \hat{\boldsymbol{\beta}} = \mathbf{V} \hat{\boldsymbol{\alpha}}.$$

Además,

$$\|\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2 = \|\mathbf{Y} - \mathbf{U} \mathbf{D} \mathbf{V}^\top \hat{\boldsymbol{\beta}}\|^2 = \|\mathbf{Z} - \mathbf{D} \hat{\boldsymbol{\alpha}}\|^2.$$

Finalmente,

$$\text{Cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = \sigma^2 (\mathbf{V} \mathbf{D}^2 \mathbf{V}^\top)^{-1} = \sigma^2 \mathbf{V} \mathbf{D}^{-2} \mathbf{V}^\top.$$

OBSERVACIÓN 2.24. Interesantemente, la SVD permite manipular problemas de rango deficiente. En efecto, cuando $\text{rg}(\mathbf{X}) < p$ podemos considerar

$$\hat{\boldsymbol{\alpha}} = \mathbf{D}^- \mathbf{Z},$$

con \mathbf{D}^- una inversa generalizada de \mathbf{D} por ejemplo

$$\mathbf{D}^- = \text{diag}(1/\delta_1, \dots, 1/\delta_r, 0, \dots, 0), \quad r = \text{rg}(\mathbf{X}),$$

y luego obtener $\hat{\boldsymbol{\beta}} = \mathbf{V} \hat{\boldsymbol{\alpha}}$.

El último procedimiento de estimación que revisaremos en esta sección corresponde al método de Gradientes Conjugados (CG), el que permite optimizar la siguiente función objetivo:

$$\phi(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}\|^2 = \frac{1}{2} (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}).$$

El algoritmo básico produce la secuencia de estimaciones,

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

con el siguiente ‘largo de paso’

$$\lambda_k = \frac{\mathbf{p}_k^\top \mathbf{g}_k}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k}, \quad \mathbf{g}_k = \mathbf{X}^\top (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}^{(k)}).$$

y actualizamos la dirección de búsqueda como:

$$\mathbf{p}_{k+1} = \mathbf{g}_{k+1} + \delta_k \mathbf{p}_k, \quad \delta_k = -\frac{\mathbf{g}_{k+1}^\top \mathbf{p}_k}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k}.$$

Para el contexto de regresión se ha sugerido modificar la versión básica del algoritmo, considerando (ver [McIntosh, 1982](#))

$$\lambda_k = \frac{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{Y}}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k},$$

y actualizar la dirección de búsqueda,

$$\mathbf{p}_{k+1} = \mathbf{g}_{k+1} + \delta_{k+1} \mathbf{p}_k, \quad \delta_{k+1} = -\frac{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{g}_k}{\mathbf{p}_k^\top \mathbf{X}^\top \mathbf{X} \mathbf{p}_k}.$$

Para hacer el proceso más simple es recomendable calcular el vector

$$\mathbf{h}_k = \mathbf{X}^\top \mathbf{X} \mathbf{p}_k,$$

lo que lleva a una implementación que solo requiere operaciones matriz-vector. Además, debemos destacar que no hace falta formar la matriz $\mathbf{X}^\top \mathbf{X}$. Lo que permite notar que los requerimientos de almacenamiento del algoritmo sólo es de $4p$ ubicaciones de memoria.

Algoritmo 3: Gradientes conjugados para regresión lineal.

Entrada : Datos \mathbf{X} y \mathbf{Y}

Parámetros: Tolerancia τ .

```

1 begin
2   Hacer  $\boldsymbol{\beta} = \mathbf{0}$ ,  $\mathbf{p} = \mathbf{g} = -\mathbf{X}^\top \mathbf{Y}$ ,  $\delta = 0$  y  $\gamma = \|\mathbf{g}\|^2$ 
3   while  $\gamma > \tau$  do
4     Calcular  $\mathbf{h} = \mathbf{X}^\top \mathbf{X} \mathbf{p}$  y  $u = \mathbf{p}^\top \mathbf{X}^\top \mathbf{X} \mathbf{p} = \mathbf{p}^\top \mathbf{h}$ 
5      $\mathbf{v} = \mathbf{g}^\top \mathbf{g}$ 
6      $\lambda = -v/u$ 
7      $\boldsymbol{\beta} = \boldsymbol{\beta} + \lambda \mathbf{p}$ 
8      $\mathbf{g} = \mathbf{g} + \lambda \mathbf{h}$ 
9      $\delta = \mathbf{g}^\top \mathbf{g} / v$ 
10     $\mathbf{p} = \mathbf{g} + \delta \mathbf{p}$ 
11  end
12  return  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}$ 
13 end
```

Los cinco procedimientos descritos en esta sección han sido implementados en la función `ols` en la biblioteca `fastmatrix` ([Osorio y Ogueda, 2021](#)) disponible para el ambiente de cálculo estadístico R ([R Core Team, 2020](#)).

2.4. Estimación bajo restricciones lineales

El objetivo de esta sección es abordar la estimación de β y σ^2 sujeto a restricciones lineales del tipo:

$$\mathbf{G}\beta = \mathbf{g}, \quad (2.8)$$

donde $\mathbf{G} \in \mathbb{R}^{q \times p}$ con $\text{rg}(\mathbf{G}) = q$ y $\mathbf{g} \in \mathbb{R}^q$. Consideraremos dos procedimientos para obtener estimadores restringidos, a saber:

- Método del modelo reducido.
- Método de multiplicadores de Lagrange.

Además, estudiaremos las propiedades estadísticas de los estimadores.

2.4.1. Método del modelo reducido. Para introducir este procedimiento, considere la siguiente partición $\mathbf{G} = (\mathbf{G}_r, \mathbf{G}_q)$ donde $\mathbf{G}_q \in \mathbb{R}^{q \times q}$ de rango q . De este modo, podemos escribir las restricciones en (2.8) como:

$$\mathbf{G}\beta = (\mathbf{G}_r, \mathbf{G}_q) \begin{pmatrix} \beta_r \\ \beta_q \end{pmatrix} = \mathbf{G}_r\beta_r + \mathbf{G}_q\beta_q = \mathbf{g},$$

como \mathbf{G}_q es no singular, tenemos

$$\beta_q = \mathbf{G}_q^{-1}(\mathbf{g} - \mathbf{G}_r\beta_r).$$

Particionando \mathbf{X} del mismo modo que $\beta = (\beta_r^\top, \beta_q^\top)^\top$, sigue que

$$\begin{aligned} \mathbf{X}\beta &= (\mathbf{X}_r, \mathbf{X}_q) \begin{pmatrix} \beta_r \\ \beta_q \end{pmatrix} = \mathbf{X}_r\beta_r + \mathbf{X}_q\beta_q \\ &= \mathbf{X}_r\beta_r + \mathbf{X}_q\mathbf{G}_q^{-1}(\mathbf{g} - \mathbf{G}_r\beta_r) \\ &= (\mathbf{X}_r - \mathbf{X}_q\mathbf{G}_q^{-1}\mathbf{G}_r)\beta_r + \mathbf{X}_q\mathbf{G}_q^{-1}\mathbf{g} \end{aligned}$$

Es decir, podemos escribir el modelo lineal

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon,$$

como:

$$\mathbf{Y} = (\mathbf{X}_r - \mathbf{X}_q\mathbf{G}_q^{-1}\mathbf{G}_r)\beta_r + \mathbf{X}_q\mathbf{G}_q^{-1}\mathbf{g} + \epsilon,$$

esto lleva al *modelo reducido*, dado por

$$\mathbf{Y}_R = \mathbf{X}_R\beta_r + \epsilon, \quad (2.9)$$

donde

$$\mathbf{Y}_R = \mathbf{Y} - \mathbf{X}_q\mathbf{G}_q^{-1}\mathbf{g}, \quad \mathbf{X}_R = \mathbf{X}_r - \mathbf{X}_q\mathbf{G}_q^{-1}\mathbf{G}_r,$$

corresponde al vector de respuesta y la matriz de diseño en el modelo reducido. La principal ventaja de este procedimiento es que permite obtener estimadores en el modelo (2.9) de forma simple. En efecto,

$$\begin{aligned} \tilde{\beta}_r &= (\mathbf{X}_R^\top \mathbf{X}_R)^{-1} \mathbf{X}_R^\top \mathbf{Y}_R, \\ s_r^2 &= \frac{1}{n-r} Q_R(\tilde{\beta}_r), \end{aligned}$$

con

$$Q_R(\tilde{\beta}_r) = \mathbf{Y}_R^\top (\mathbf{I} - \mathbf{X}_R(\mathbf{X}_R^\top \mathbf{X}_R)^{-1} \mathbf{X}_R^\top) \mathbf{Y}_R$$

Por otro lado, el vector de coeficientes estimados bajo las restricciones lineales en (2.8)

$$\begin{aligned}\tilde{\beta} &= \begin{pmatrix} \tilde{\beta}_r \\ \tilde{\beta}_q \end{pmatrix} = \begin{pmatrix} \tilde{\beta}_r \\ \mathbf{G}_q^{-1}(\mathbf{g} - \mathbf{G}_r \tilde{\beta}_r) \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{0} \\ \mathbf{G}_q^{-1} \mathbf{g} \end{pmatrix} + \begin{pmatrix} \mathbf{I} \\ -\mathbf{G}_q^{-1} \mathbf{G}_r \end{pmatrix} \tilde{\beta}_r,\end{aligned}\quad (2.10)$$

Mientras que el vector de ‘residuos’ en el modelo reducido adopta la forma,

$$\begin{aligned}\mathbf{Y}_R - \mathbf{X}_R \tilde{\beta}_r &= \mathbf{Y} - \mathbf{X}_q \mathbf{G}_q^{-1} \mathbf{g} - (\mathbf{X}_r - \mathbf{X}_q \mathbf{G}_q^{-1} \mathbf{g}) \tilde{\beta}_r \\ &= \mathbf{Y} - (\mathbf{X}_r, \mathbf{X}_q) \begin{pmatrix} \tilde{\beta}_r \\ \tilde{\beta}_q \end{pmatrix} = \mathbf{Y} - \mathbf{X} \tilde{\beta},\end{aligned}$$

de este modo

$$Q_R(\tilde{\beta}_r) = \|\mathbf{Y}_R - \mathbf{X}_R \tilde{\beta}_r\|^2 = \|\mathbf{Y} - \mathbf{X} \tilde{\beta}\|^2 = Q(\tilde{\beta}). \quad (2.11)$$

Las expresiones anteriores permiten establecer el siguiente resultado.

RESULTADO 2.25. *Para el modelo lineal $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ sujeto a las restricciones $\mathbf{G}\beta = \mathbf{g}$ con $\epsilon \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. El MLE restringido de β es dado por (2.10) y tenemos que*

$$\tilde{\beta} \sim \mathbf{N}_p(\beta, \text{Cov}(\tilde{\beta})),$$

donde

$$\text{Cov}(\tilde{\beta}) = \sigma^2 \begin{pmatrix} \mathbf{I} \\ -\mathbf{G}_q^{-1} \mathbf{G}_r \end{pmatrix} (\mathbf{X}_R^\top \mathbf{X}_R)^{-1} \begin{pmatrix} \mathbf{I} \\ -\mathbf{G}_q^{-1} \mathbf{G}_r \end{pmatrix}^\top.$$

Mientras que

$$s_r^2 = \frac{1}{n-r} Q_R(\tilde{\beta}_r) = \frac{1}{n-r} Q(\tilde{\beta}),$$

donde

$$\frac{(n-r)s_r^2}{\sigma^2} \sim \chi^2(n-r),$$

y $\tilde{\beta}$, $Q(\tilde{\beta})$ son independientes.

DEMOSTRACIÓN. Sabemos que

$$\begin{aligned}\tilde{\beta}_r &\sim \mathbf{N}_r(\beta_r, \sigma^2 (\mathbf{X}_R^\top \mathbf{X}_R)^{-1}), \\ \frac{(n-r)s_r^2}{\sigma^2} &= \frac{Q_R(\tilde{\beta}_r)}{\sigma^2} = \frac{\mathbf{Y}_R^\top (\mathbf{I} - \mathbf{H}_R) \mathbf{Y}_R}{\sigma^2} \sim \chi^2(n-r).\end{aligned}$$

Así, por (2.10), tenemos

$$\begin{aligned}\mathbb{E}(\tilde{\beta}) &= \begin{pmatrix} \mathbf{0} \\ \mathbf{G}_q^{-1} \mathbf{g} \end{pmatrix} + \begin{pmatrix} \mathbf{I}_r \\ -\mathbf{G}_q^{-1} \mathbf{G}_r \end{pmatrix} \mathbb{E}(\tilde{\beta}_r) \\ &= \begin{pmatrix} \beta_r \\ \mathbf{G}_q^{-1}(\mathbf{g} - \mathbf{G}_r \beta_r) \end{pmatrix} = \begin{pmatrix} \beta_r \\ \beta_q \end{pmatrix} = \beta\end{aligned}$$

Además,

$$\begin{aligned}\text{Cov}(\tilde{\beta}) &= \begin{pmatrix} \mathbf{I}_r \\ -\mathbf{G}_q^{-1}\mathbf{G}_r \end{pmatrix} \text{Cov}(\tilde{\beta}_r) \begin{pmatrix} \mathbf{I}_r \\ -\mathbf{G}_q^{-1}\mathbf{G}_r \end{pmatrix}^\top \\ &= \sigma^2 \begin{pmatrix} \mathbf{I}_r \\ -\mathbf{G}_q^{-1}\mathbf{G}_r \end{pmatrix} (\mathbf{X}_R^\top \mathbf{X}_R)^{-1} \begin{pmatrix} \mathbf{I}_r \\ -\mathbf{G}_q^{-1}\mathbf{G}_r \end{pmatrix}^\top.\end{aligned}$$

como $\tilde{\beta}$ es una función lineal de $\tilde{\beta}$ la normalidad sigue. La independencia entre $\tilde{\beta}$ y $Q(\tilde{\beta})$ es consecuencia del Resultado 2.7. \square

2.4.2. Método de multiplicadores de Lagrange. Considere

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon, \quad \epsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}).$$

La función Langrangiana asociada a las restricciones lineales $\mathbf{G}\beta = \mathbf{g}$ es dada por:

$$F(\theta, \lambda) = \ell(\theta) + \frac{1}{\sigma^2} \lambda^\top (\mathbf{G}\beta - \mathbf{g}),$$

con $\theta = (\beta^\top, \sigma^2)^\top$. De este modo,

$$\begin{aligned}\frac{\partial F(\theta, \lambda)}{\partial \beta} &= \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta) + \frac{1}{\sigma^2} \mathbf{G}^\top \lambda \\ \frac{\partial F(\theta, \lambda)}{\partial \sigma^2} &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \{Q(\beta) - 2\lambda^\top (\mathbf{G}\beta - \mathbf{g})\} \\ \frac{\partial F(\theta, \lambda)}{\partial \lambda} &= \mathbf{G}\beta - \mathbf{g}\end{aligned}$$

Desde la condición de primer orden, obtenemos las ecuaciones de estimación,

$$\begin{aligned}\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta) + \mathbf{G}^\top \lambda &= \mathbf{0}, \\ n\sigma^2 - \{Q(\beta) - 2\lambda^\top (\mathbf{G}\beta - \mathbf{g})\} &= 0, \\ \mathbf{G}\beta &= \mathbf{g},\end{aligned}$$

es decir,

$$\mathbf{X}^\top \mathbf{X}\beta = \mathbf{X}^\top \mathbf{Y} + \mathbf{G}^\top \lambda, \quad (2.12)$$

$$\sigma^2 = \frac{1}{n} Q(\beta) \quad (2.13)$$

$$\mathbf{G}\beta = \mathbf{g}, \quad (2.14)$$

Resolviendo la Ecuación (2.12) con relación a β obtenemos

$$\tilde{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{Y} + \mathbf{G}^\top \lambda)$$

Substituyendo este resultado en (2.14) y resolviendo para λ , sigue que

$$\mathbf{G}\tilde{\beta} = \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{Y} + \mathbf{G}^\top \lambda) = \mathbf{g},$$

es decir,

$$\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} + \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top \lambda = \mathbf{g},$$

por tanto,

$$\tilde{\lambda} = (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{g} - \mathbf{G}\hat{\beta}).$$

Reemplazando este resultado en $\tilde{\beta}$ resulta

$$\begin{aligned}\tilde{\beta} &= (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{Y} + \mathbf{G}^\top \tilde{\lambda}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{g} - \mathbf{G}\hat{\beta}) \\ &= \hat{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{g} - \mathbf{G}\hat{\beta}).\end{aligned}$$

Que puede ser reorganizado como:

$$\tilde{\beta} = \mathbf{A}\hat{\beta} + \mathbf{B}\mathbf{g} = \hat{\beta} - \mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g}),$$

donde $\hat{\beta}$ corresponde al MLE no restringido para β , con

$$\mathbf{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} \quad (2.15)$$

$$\mathbf{A} = \mathbf{I} - \mathbf{B}\mathbf{G} \quad (2.16)$$

y el estimador insesgado para σ^2 es dado por

$$s_r^2 = \frac{1}{n-r} Q(\tilde{\beta}).$$

Para estudiar las propiedades de este MLE restringido, considere primeramente el siguiente lema.

LEMA 2.26. *La matriz \mathbf{A} definida en (2.16) tiene las siguientes propiedades:*

- (i) \mathbf{A} es idempotente con $\text{rg}(\mathbf{A}) = r$.
- (ii) $\mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ es idempotente y simétrica con rango r .
- (iii) $\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top = \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top$.

DEMOSTRACIÓN. Para mostrar que $\mathbf{A} = \mathbf{I} - \mathbf{B}\mathbf{G}$ es idempotente, basta mostrar que $\mathbf{B}\mathbf{G}$ es idempotente. En efecto,

$$\mathbf{G}\mathbf{B}\mathbf{G} = \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} \mathbf{G} = \mathbf{G},$$

de ahí que $\mathbf{B}\mathbf{G}$ es idempotente. Además,

$$\text{rg}(\mathbf{B}\mathbf{G}) = \text{tr}(\mathbf{B}\mathbf{G}) = \text{tr}(\mathbf{G}\mathbf{B}) = q,$$

así $\text{rg}(\mathbf{I} - \mathbf{B}\mathbf{G}) = \text{tr}(\mathbf{A}) = p - q = r$, lo que muestra la parte (i).

Para notar la parte (ii), sea $\mathbf{C} = \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$, luego

$$\mathbf{C}^2 = \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{X}\mathbf{A}^2(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top,$$

como \mathbf{A} es idempotente, sigue que $\mathbf{C}^2 = \mathbf{C}$. Esto permite escribir

$$\text{rg}(\mathbf{C}) = \text{tr}(\mathbf{C}) = \text{tr}(\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}) = \text{tr}(\mathbf{A}) = r.$$

Por otro lado,

$$\mathbf{C}^\top = (\mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top \mathbf{X}^\top.$$

Tenemos que,

$$\mathbf{A}^\top = \mathbf{I} - \mathbf{G}^\top \mathbf{B}^\top = \mathbf{I} - \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1},$$

luego premultiplicando por $(\mathbf{X}^\top \mathbf{X})^{-1}$ y factorizando lleva a,

$$\begin{aligned}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top &= (\mathbf{X}^\top \mathbf{X})^{-1} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{I} - \mathbf{B}\mathbf{G})(\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1},\end{aligned} \quad (2.17)$$

así

$$\mathbf{C}^\top = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top \mathbf{X} = \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X} = \mathbf{C}.$$

Finalmente, Ecuación (2.17) permite notar la primera igualdad de la parte (iii). Ahora, por (2.17), sigue que

$$\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top = \mathbf{A}^2(\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1},$$

pues \mathbf{A} es idempotente y esto termina la prueba. \square

Esto nos habilita para establecer el siguiente resultado,

RESULTADO 2.27. *Para el modelo lineal*

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}).$$

El MLE de $\boldsymbol{\beta}$ bajo las restricciones lineales $\mathbf{G}\boldsymbol{\beta} = \mathbf{g}$, es dado por

$$\tilde{\boldsymbol{\beta}} = \mathbf{A}\hat{\boldsymbol{\beta}} + \mathbf{B}\mathbf{g},$$

con distribución

$$\tilde{\boldsymbol{\beta}} \sim \mathbf{N}_p(\boldsymbol{\beta}, \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top).$$

Mientras que el estimador insesgado de σ^2 es

$$s_r^2 = \frac{1}{n-r} Q(\tilde{\boldsymbol{\beta}}),$$

donde

$$\frac{Q(\tilde{\boldsymbol{\beta}})}{\sigma^2} \sim \chi^2(n-r),$$

y $\tilde{\boldsymbol{\beta}}$ es independiente de $Q(\tilde{\boldsymbol{\beta}})$.

DEMOSTRACIÓN. La normalidad de $\tilde{\boldsymbol{\beta}}$ sigue desde la linealidad con relación a $\hat{\boldsymbol{\beta}}$. Ahora,

$$\mathbf{E}(\tilde{\boldsymbol{\beta}}) = \mathbf{A} \mathbf{E}(\hat{\boldsymbol{\beta}}) + \mathbf{B}\mathbf{g} = (\mathbf{I} - \mathbf{B}\mathbf{G})\boldsymbol{\beta} + \mathbf{B}\mathbf{g} = \boldsymbol{\beta} - \mathbf{B}(\mathbf{G}\boldsymbol{\beta} - \mathbf{g}) = \boldsymbol{\beta},$$

y

$$\text{Cov}(\tilde{\boldsymbol{\beta}}) = \mathbf{A} \text{Cov}(\hat{\boldsymbol{\beta}}) \mathbf{A}^\top = \sigma^2 \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top.$$

Notando que

$$\begin{aligned} \tilde{\boldsymbol{\beta}} &= \mathbf{A}\hat{\boldsymbol{\beta}} + \mathbf{B}\mathbf{g} = \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} + \mathbf{B}\mathbf{g} \\ &= \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) + \mathbf{B}\mathbf{g} \\ &= \mathbf{A}\boldsymbol{\beta} + \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} + \mathbf{B}\mathbf{g} \\ &= \boldsymbol{\beta} + \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{aligned}$$

lo que permite escribir

$$\begin{aligned} \mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}} &= \mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ &= \boldsymbol{\epsilon} - \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ &= (\mathbf{I} - \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \boldsymbol{\epsilon}, \end{aligned}$$

por la parte (ii) del Lema 2.26, sigue que

$$\frac{Q(\tilde{\boldsymbol{\beta}})}{\sigma^2} = \frac{\boldsymbol{\epsilon}^\top (\mathbf{I} - \mathbf{X}\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \boldsymbol{\epsilon}}{\sigma^2} \sim \chi^2(n-r).$$

Para notar la independencia entre $\tilde{\beta}$ y $Q(\tilde{\beta})$ debemos tener:⁶

$$\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{X} \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \mathbf{0}.$$

En efecto,

$$\begin{aligned} & \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{X} \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\ &= \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} - \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top. \end{aligned}$$

Notando que $\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top$, obtenemos

$$\begin{aligned} & \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top \mathbf{X}^\top \\ &= \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top \mathbf{X}^\top = \mathbf{A}^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top, \end{aligned}$$

y esto concluye la demostración. \square

2.5. Test de hipótesis lineales

El objetivo de esta sección recae en desarrollar el test de razón de verosimilitudes para probar hipótesis lineales de la forma:

$$H_0 : \mathbf{G}\beta = \mathbf{g}, \quad \text{versus} \quad H_1 : \mathbf{G}\beta \neq \mathbf{g} \quad (2.18)$$

donde \mathbf{G} es una matriz de contrastes de orden $q \times p$ con $\text{rg}(\mathbf{G}) = q$ y $\mathbf{g} \in \mathbb{R}^q$.

OBSERVACIÓN 2.28. H_0 es expresada como un sistema de ecuaciones mientras que H_1 indica que al menos una ecuación no se satisface.

Para abordar hipótesis lineales, usaremos el principio de verosimilitud. Es decir, consideraremos el estadístico

$$\Lambda = \frac{\max_{\mathbf{G}\beta=\mathbf{g}} L(\beta, \sigma^2)}{\max_{\Theta} L(\beta, \sigma^2)} = \frac{L(\tilde{\beta}, \tilde{\sigma}^2)}{L(\hat{\beta}, \hat{\sigma}^2)}.$$

Asumiendo $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{I})$, tenemos la función de verosimilitud

$$L(\beta, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) \right\}.$$

Sabemos que

$$\begin{aligned} \max_{\Theta} L(\beta, \sigma^2) &= L(\hat{\beta}, \hat{\sigma}^2) \\ &= (2\pi\hat{\sigma}^2)^{-n/2} \exp \left\{ -\frac{1}{2\hat{\sigma}^2} \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 \right\} \\ &= \{2\pi Q(\hat{\beta})/n\}^{-n/2} \exp \left\{ -\frac{n}{2Q(\hat{\beta})} \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 \right\} \\ &= \{2\pi Q(\hat{\beta})/n\}^{-n/2} \exp(-n/2). \end{aligned}$$

Mientras que bajo $H_0 : \mathbf{G}\beta = \mathbf{g}$, tenemos⁷

$$\begin{aligned} \max_{\mathbf{G}\beta=\mathbf{g}} L(\beta, \sigma^2) &= L(\tilde{\beta}, \tilde{\sigma}^2) = (2\pi\tilde{\sigma}^2)^{-n/2} \exp \left\{ -\frac{1}{2\tilde{\sigma}^2} \|\mathbf{Y} - \mathbf{X}\tilde{\beta}\|^2 \right\} \\ &= \{2\pi Q(\tilde{\beta})/n\}^{-n/2} \exp(-n/2). \end{aligned}$$

⁶Lo que es consecuencia de escribir $\tilde{\beta}$ y $Q(\tilde{\beta})$ en términos de ϵ .

⁷Con espacio paramétrico nulo, $\Theta_0 = \{\theta = (\beta^\top, \sigma^2)^\top : \mathbf{G}\beta = \mathbf{g}\}$.

De este modo, el estadístico de razón de verosimilitudes

$$\begin{aligned}\Lambda &= \frac{L(\tilde{\beta}, \tilde{\sigma}^2)}{L(\hat{\beta}, \hat{\sigma}^2)} = \frac{\{2\pi Q(\tilde{\beta})/n\}^{-n/2} \exp(-n/2)}{\{2\pi Q(\hat{\beta})/n\}^{-n/2} \exp(-n/2)} \\ &= \left\{ \frac{Q(\hat{\beta})}{Q(\tilde{\beta})} \right\}^{n/2},\end{aligned}$$

y de acuerdo con el principio de verosimilitud rechazamos $H_0 : \mathbf{G}\beta = \mathbf{g}$ si Λ es pequeño.

Alternativamente, podemos considerar

$$\Lambda^{2/n} = \frac{Q(\hat{\beta})}{Q(\tilde{\beta})}.$$

Recuerde que $\tilde{\beta} = \hat{\beta} - \mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g})$, así

$$\mathbf{Y} - \mathbf{X}\tilde{\beta} = \mathbf{Y} - \mathbf{X}(\hat{\beta} - \mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g})) = \mathbf{Y} - \mathbf{X}\hat{\beta} + \mathbf{X}\mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g}).$$

Además,

$$\begin{aligned}Q(\tilde{\beta}) &= \|\mathbf{Y} - \mathbf{X}\tilde{\beta}\|^2 = \|\mathbf{Y} - \hat{\beta} - \mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g})\|^2 \\ &= \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 + \|\mathbf{X}\mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g})\|^2 \\ &\quad + (\mathbf{Y} - \mathbf{X}\hat{\beta})^\top \mathbf{X}\mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g}) + (\mathbf{G}\hat{\beta} - \mathbf{g})^\top \mathbf{B}^\top \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}).\end{aligned}$$

Sin embargo,

$$\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}) = \mathbf{X}^\top (\mathbf{I} - \mathbf{H})\mathbf{Y} = \mathbf{0}$$

lo que nos lleva a:

$$\begin{aligned}Q(\tilde{\beta}) &= Q(\hat{\beta}) + (\mathbf{G}\hat{\beta} - \mathbf{g})^\top \mathbf{B}^\top \mathbf{X}^\top \mathbf{X}\mathbf{B}(\mathbf{G}\hat{\beta} - \mathbf{g}) \\ &\geq Q(\hat{\beta})\end{aligned}$$

OBSERVACIÓN 2.29. Es decir,

$$0 \leq \frac{Q(\hat{\beta})}{Q(\tilde{\beta})} \leq 1,$$

por tanto, $\Lambda^{2/n} \in [0, 1]$.

Además, recordando que

$$\mathbf{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1},$$

obtenemos

$$\mathbf{B}^\top \mathbf{X}^\top \mathbf{X}\mathbf{B} = (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1}.$$

De este modo,

$$Q(\tilde{\beta}) - Q(\hat{\beta}) = (\mathbf{G}\hat{\beta} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\beta} - \mathbf{g}).$$

Así

$$\frac{Q(\tilde{\beta}) - Q(\hat{\beta})}{Q(\hat{\beta})} = \Lambda^{-2/n} - 1,$$

es decir, valores pequeños de Λ (en cuyo caso rechazamos H_0) implican valores grandes de la razón anterior.

Ahora considere

$$\begin{aligned} E(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}) &= \mathbf{G} E(\hat{\boldsymbol{\beta}}) - \mathbf{g} = \mathbf{G}\boldsymbol{\beta} - \mathbf{g} \\ \text{Cov}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}) &= \mathbf{G} \text{Cov}(\hat{\boldsymbol{\beta}}) \mathbf{G}^\top = \sigma^2 \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top. \end{aligned}$$

Note que

$$(\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} = \left\{ \frac{1}{\sigma^2} \text{Cov}(\mathbf{G}\hat{\boldsymbol{\beta}}) \right\}^{-1}.$$

De esta forma

$$\frac{(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})}{\sigma^2} \sim \chi^2(q; \delta),$$

pues $\sigma^{-2}(\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} \text{Cov}(\mathbf{G}\hat{\boldsymbol{\beta}}) = \mathbf{I}$, es matriz idempotente, y

$$\delta = \frac{1}{2\sigma^2} (\mathbf{G}\boldsymbol{\beta} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\boldsymbol{\beta} - \mathbf{g})$$

Por otro lado, sabemos que

$$\frac{Q(\hat{\boldsymbol{\beta}})}{\sigma^2} \sim \chi^2(n - p; 0)$$

Para notar la independencia, considere $\boldsymbol{\beta}^*$ cualquier vector que satisface la condición $\mathbf{G}\boldsymbol{\beta}^* = \mathbf{g}$. Entonces,

$$\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I} - \mathbf{H})\mathbf{Y} = (\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*),$$

pues $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ y

$$\begin{aligned} \mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g} &= \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - \mathbf{G}\boldsymbol{\beta}^* = \mathbf{G}\{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - \boldsymbol{\beta}^*\} \\ &= \mathbf{G}\{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}^*\} \\ &= \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*). \end{aligned}$$

En nuestro caso,

$$\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^* \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\beta}^*, \sigma^2 \mathbf{I}) \stackrel{d}{=} \mathbf{N}_n(\mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\beta}^*), \sigma^2 \mathbf{I}).$$

De este modo,

$$Q(\hat{\boldsymbol{\beta}}) = \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*)^\top (\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*),$$

mientras que

$$\begin{aligned} Q(\tilde{\boldsymbol{\beta}}) - Q(\hat{\boldsymbol{\beta}}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*)^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} \\ &\quad \times \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*). \end{aligned}$$

Como

$$(\mathbf{I} - \mathbf{H})\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} \mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{0},$$

sigue la independencia y permite construir la estadística

$$F = \frac{\{Q(\tilde{\boldsymbol{\beta}}) - Q(\hat{\boldsymbol{\beta}})\}/q}{Q(\hat{\boldsymbol{\beta}})/(n - p)} \sim F(q, n - p; \delta).$$

Esto lleva al siguiente resultado.

RESULTADO 2.30. Para el modelo lineal $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ con los supuestos A1 a A4*. Un test de tamaño α para probar

$$H_0 : \mathbf{G}\boldsymbol{\beta} = \mathbf{g}, \quad \text{versus} \quad H_1 : \mathbf{G}\boldsymbol{\beta} \neq \mathbf{g},$$

es dado por, rechazar H_0 cuando:

$$F = \frac{Q(\tilde{\boldsymbol{\beta}}) - Q(\hat{\boldsymbol{\beta}})}{qs^2} \geq F_{1-\alpha}(q, n-p; 0).$$

DEMOSTRACIÓN. Bajo $H_0 : \mathbf{G}\boldsymbol{\beta} = \mathbf{g}$, tenemos $\delta = 0$, de ahí que $F \sim F(q, n-p; 0)$, lo que lleva al resultado deseado. \square

OBSERVACIÓN 2.31. Note que podemos escribir el estadístico F de varias formas equivalentes, a saber:

$$\begin{aligned} F &= \left(\frac{n-p}{q} \right) \frac{Q(\tilde{\boldsymbol{\beta}}) - Q(\hat{\boldsymbol{\beta}})}{Q(\hat{\boldsymbol{\beta}})} \\ &= \frac{Q(\tilde{\boldsymbol{\beta}}) - Q(\hat{\boldsymbol{\beta}})}{qs^2} \\ &= \frac{(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})}{qs^2} \sim F(q, n-p, \delta). \end{aligned}$$

Hemos notado la relación que existe entre el test de razón de verosimilitudes con el estadístico F para probar hipótesis lineales en el modelo de regresión lineal

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}).$$

A continuación exploramos la relación entre el estadístico F con los *test de Wald, score y gradiente* para hipótesis lineales, del tipo

$$H_0 : \mathbf{G}\boldsymbol{\beta} = \mathbf{g}, \quad \text{versus} \quad H_1 : \mathbf{G}\boldsymbol{\beta} \neq \mathbf{g}.$$

Primeramente, note que la matriz de información de Fisher para $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$, adopta la forma:

$$\mathcal{F}(\boldsymbol{\theta}) = \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \frac{n}{2\sigma^2} \end{pmatrix}.$$

El estadístico de Wald para hipótesis lineales de la forma $H_0 : \mathbf{G}\boldsymbol{\beta} = \mathbf{g}$, es dado por

$$\begin{aligned} W &= n(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top (\mathbf{G}\{\mathcal{F}(\hat{\boldsymbol{\beta}})\}^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}) \\ &= \frac{n(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})}{\hat{\sigma}^2}. \end{aligned}$$

Mientras que el test score es dado por

$$\begin{aligned} R &= \frac{1}{n} \mathbf{U}^\top(\tilde{\boldsymbol{\beta}}) \{\mathcal{F}(\tilde{\boldsymbol{\beta}})\}^{-1} \mathbf{U}(\tilde{\boldsymbol{\beta}}) \\ &= \frac{1}{n\hat{\sigma}^2} (\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}}). \end{aligned}$$

Como,

$$\begin{aligned} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) &= \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})) \\ &= \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + \mathbf{X}^\top \mathbf{X}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}) \\ &= \mathbf{X}^\top \mathbf{X}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}). \end{aligned}$$

Sigue que,

$$\begin{aligned}
& (\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) \\
&= (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top \mathbf{B}^\top \mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{B}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}) \\
&= (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top \mathbf{B}^\top \mathbf{X}^\top \mathbf{X} \mathbf{B}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}) \\
&= (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}).
\end{aligned}$$

Finalmente,

$$R = \frac{(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})}{n\tilde{\sigma}^2}$$

Por otro lado, el estadístico gradiente es dado por:

$$T = \mathbf{U}^\top (\tilde{\boldsymbol{\beta}})(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}) = (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top \mathbf{B}^\top \mathbf{X}^\top \mathbf{X}(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})$$

Sabemos que $\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - \mathbf{B}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})$, esto lleva a

$$\begin{aligned}
T &= (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top \mathbf{B}^\top \mathbf{X}^\top \mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}} + \mathbf{B}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})) \\
&= (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top \mathbf{B}^\top \mathbf{X}^\top \mathbf{X} \mathbf{B}(\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g}) \\
&= (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\mathbf{G}\hat{\boldsymbol{\beta}} - \mathbf{g})
\end{aligned}$$

OBSERVACIÓN 2.32. Es decir, podemos escribir:

$$W = \frac{qn}{n-p} F, \quad R = \left\{ 1 + \left(\frac{n-p}{q} \right) F^{-1} \right\}^{-1}, \quad T = \frac{qs^2}{n-p} F.$$

Lo que permite notar que basta usar el estadístico F para probar hipótesis lineales en el modelo de regresión lineal.

2.6. Regiones de confianza

Una región de confianza del $100(1 - \alpha)\%$ para $\boldsymbol{\gamma}$ se define como la región en el espacio paramétrico, digamos $RC(\boldsymbol{\gamma})$ con la propiedad

$$P(\boldsymbol{\gamma} \in RC(\boldsymbol{\gamma})) = 1 - \alpha,$$

donde $\boldsymbol{\gamma}$ es el verdadero vector de parámetros.

Sea $\boldsymbol{\gamma} = \mathbf{G}\boldsymbol{\beta}$ un vector de parámetros q -dimensional. Sabemos que

$$\frac{(\boldsymbol{\gamma} - \hat{\boldsymbol{\gamma}})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\boldsymbol{\gamma} - \hat{\boldsymbol{\gamma}})}{qs^2} \sim F(q, n-p),$$

luego, puede ser usada como una región de confianza. Es decir,

$$CR(\boldsymbol{\gamma}) = \{ \boldsymbol{\gamma} : (\boldsymbol{\gamma} - \hat{\boldsymbol{\gamma}})^\top (\mathbf{G}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{G}^\top)^{-1} (\boldsymbol{\gamma} - \hat{\boldsymbol{\gamma}}) \leq qs^2 F_{1-\alpha}(q, n-p) \},$$

donde $F_{1-\alpha}(q, n-p)$ es el valor cuantil $1 - \alpha$ de la distribución F con q y $n-p$ grados de libertad. De este modo, para $\mathbf{G} = \mathbf{I}$ obtenemos el siguiente caso particular

$$CR(\boldsymbol{\beta}) = \{ \boldsymbol{\beta} : (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \leq ps^2 F_{1-\alpha}(p, n-p) \}.$$

Notando que

$$\frac{\mathbf{a}^\top \widehat{\boldsymbol{\beta}} - \mathbf{a}^\top \boldsymbol{\beta}}{s \sqrt{\mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}} \sim t(n-p), \quad \mathbf{a} \in \mathbb{R}^p,$$

corresponde a una cantidad pivotal, podemos escribir un intervalo de confianza para la combinación lineal $\tau = \mathbf{a}^\top \boldsymbol{\beta}$. Esto es,

$$\tau \in \left[\mathbf{a}^\top \widehat{\boldsymbol{\beta}} \mp t_{1-\alpha/2}(n-p) s \sqrt{\mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}} \right],$$

donde $t_{1-\alpha/2}(n-p)$ denota el valor cuantil $1 - \alpha/2$ de la distribución t de Student con $n-p$ grados de libertad. También podemos escribir de forma equivalente:

$$CR(\mathbf{a}^\top \boldsymbol{\beta}) = \{ \tau : (\tau - \mathbf{a}^\top \widehat{\boldsymbol{\beta}}) \leq s^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} F_{1-\alpha}(1, n-p) \}.$$

En particular, para β_j ($j = 1, \dots, p$) basta escoger $\mathbf{a} = \mathbf{e}_j$ el j -ésimo vector unidad. Así,

$$CI(\beta_j) = [\widehat{\beta}_j \mp t_{1-\alpha/2}(n-p) s \sqrt{c_{jj}}],$$

donde c_{jj} representa el j -ésimo elemento de la diagonal de $\mathbf{C} = (\mathbf{X}^\top \mathbf{X})^{-1}$.

Ejercicios

2.1 Sean Y_1, \dots, Y_n variables aleatorias independientes con $Y_i \sim \mathbf{N}(\alpha + \theta z_i, \sigma^2)$, $i = 1, \dots, n$, donde $\{z_i\}$ son constantes conocidas, tales que $\sum_{i=1}^n z_i = 0$. Obtenga el estimador ML de $\boldsymbol{\beta} = (\alpha, \theta)^\top$ y determine su matriz de covarianza. ¿Son $\widehat{\alpha}$ y $\widehat{\theta}$ independientes?

2.2 Sea Y_{ij} , para $i = 1, 2, 3$ y $j = 1, \dots, m$ variables aleatorias independientes con distribución normal, tales que $\mathbf{E}(Y_{ij}) = \mu_{ij}$, $\text{var}(Y_{ij}) = \sigma^2$, y

$$\mu_{1j} = \tau, \quad \mu_{2j} = \tau + \theta, \quad \mu_{3j} = \tau - \theta.$$

- Determine la matriz de diseño \mathbf{X} .
- Obtener el estimador ML de $(\tau, \theta)^\top$ y $\text{var}(\widehat{\theta})$.
- Derive el estadístico F para probar la hipótesis $H_0 : \theta = 0$.

2.3 Sea $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$ y $z_{ij} = x_{ij} - \bar{x}_j$, para todo $i = 1, \dots, n$; $j = 1, \dots, k$ y considere $\mathbf{Z} = (z_{ij})$, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)^\top$. En cuyo caso, tenemos el *modelo centrado*:

$$\mathbf{Y} = \alpha \mathbf{1} + \mathbf{Z} \boldsymbol{\beta} + \boldsymbol{\epsilon} = (\mathbf{1}, \mathbf{Z}) \begin{pmatrix} \alpha \\ \boldsymbol{\beta} \end{pmatrix} + \boldsymbol{\epsilon},$$

donde $\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$ y $\text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$. Muestre que los BLUE de α y $\boldsymbol{\beta}$ son independientes.

2.4 Considere las regresiones de Y sobre x para los datos a continuación, especificadas por:

$$\mathcal{M}_1 : \mathbf{E}(Y) = \beta_0 x \quad \text{y} \quad \mathcal{M}_2 : \mathbf{E}(Y) = \beta_1 x + \beta_2 x^2.$$

Obtenga $\widehat{\beta}_0$, $\widehat{\beta}_1$ y $\widehat{\beta}_2$. ¿Cuál de esos modelos es preferido?

Y	5	7	7	10	16	20
x	1	2	3	4	5	6

2.5 Considere las rectas de regresión:

$$\mathcal{R}_1 : Y_{1i} = \alpha_1 + \beta_1 x_{1i} + \epsilon_{1i} \quad \text{y}$$

$$\mathcal{R}_2 : Y_{2i} = \alpha_2 + \beta_2 x_{2i} + \epsilon_{2i},$$

para $i = 1, \dots, n$, donde los errores $\{\epsilon_{1i}\}$ y $\{\epsilon_{2i}\}$ son variables aleatorias iid con media cero y varianza común σ^2 . Obtenga el estadístico F para probar la hipótesis de que las rectas de regresión \mathcal{R}_1 y \mathcal{R}_2 son paralelas.

2.6 Considere el modelo

$$Y_{ij} = \theta_i x_j + \epsilon_{ij}, \quad i = 1, 2; j = 1, \dots, T,$$

con $\{\epsilon_{ij}\}$ variables aleatorias independientes $N(0, \sigma^2)$ y $\{x_i\}$ constantes conocidas. Obtenga el estadístico F para probar $H_0 : \theta_1 = \theta_2$.

Chequeo del Modelo y Alternativas a Mínimos Cuadrados

En este capítulo se describe la inferencia en modelos lineales. Primeramente introducimos algunas definiciones y supuestos en los que se basan los modelos de regresión.

3.1. Colinealidad

Considere el modelo

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

donde $\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$ y $\text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$ con $\mathbf{X} \in \mathbb{R}^{n \times p}$ tal que $\text{rg}(\mathbf{X}) = p$. Es bien conocido que cuando \mathbf{X} es mal condicionada, el sistema de ecuaciones

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{Y}, \quad (3.1)$$

puede ser muy inestable.

Debemos destacar que, cuando la matriz de diseño es de rango (columna) deficiente, entonces podemos considerar alguna inversa generalizada para obtener una solución del problema en (3.1). En esta sección, nos enfocaremos en el problema en que $\text{rg}(\mathbf{X}) = p$, y sin embargo tenemos que existe \mathbf{a} tal que $\mathbf{X}\mathbf{a} \approx \mathbf{0}$.

OBSERVACIÓN 3.1. Este es un problema numérico que puede tener consecuencias inferenciales importantes, por ejemplo:

- (a) Típicamente los coeficientes estimados $\hat{\boldsymbol{\beta}}$ tendrán varianzas “grandes”.
- (b) Test estadísticos presentarán bajo poder y los intervalos de confianza serán muy amplios.
- (c) Signos de algunos coeficientes son “incorrectos” (basados en conocimiento previo).
- (d) Resultados cambian bruscamente con la eliminación de una o varias columnas de \mathbf{X} .

Algunas herramientas para el diagnóstico de colinealidad, son:

- (a) Examinar la *matriz de correlación* entre los regresores y la respuesta, esto es:

$$\begin{pmatrix} \mathbf{R}_{XX} & \mathbf{R}_{XY} \\ & 1 \end{pmatrix},$$

correlaciones altas entre dos variables regresoras pueden indicar un posible problema de colinealidad.

- (b) *Factores de inflación de varianza*: Suponga que los datos han sido centrados y escalados, entonces

$$\mathbf{R}^{-1} = (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}, \quad \widetilde{\mathbf{X}} = (x_{ij} - \bar{x}_j),$$

y los elementos diagonales de \mathbf{R}^{-1} son llamados factores de inflación de varianza VIF_j . Se puede mostrar que

$$\text{VIF}_j = \frac{1}{1 - R_j^2},$$

donde R_j^2 es el coeficiente de correlación múltiple de \mathbf{X}_j “regresado” sobre el resto de variables explicativas y de ahí que un VIF_j “alto” indica R_j^2 cercano a 1 y por tanto presencia de colinealidad.

- (c) Examinar los valores/vectores propios (o *componentes principales*) de la matriz de correlación \mathbf{R} .
- (d) *Número condición*: Desde la SVD de \mathbf{X} podemos escribir

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top,$$

donde $\mathbf{U} \in \mathbb{R}^{n \times p}$, $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_p$, $\mathbf{D} = \text{diag}(\delta_1, \dots, \delta_p)$ y $\mathbf{V} \in \mathcal{O}_p$. La detección de colinealidad puede ser llevada a cabo usando

$$\kappa(\mathbf{X}) = \|\mathbf{X}\| \|\mathbf{X}^+\| = \frac{\delta_1}{\delta_p},$$

y $\kappa(\mathbf{X})$ “grande” ($\kappa > 30$) es un indicador de colinealidad.

Note que, el caso de *deficiencia de rango* puede ser manipulado sin problemas usando SVD. En efecto,

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top = \mathbf{U} \begin{pmatrix} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V}^\top$$

donde $\mathbf{D}_1 \in \mathbb{R}^{r \times r}$, $\text{rg}(\mathbf{X}) = r < p$. De este modo

$$\mathbf{X} \mathbf{V} = \mathbf{U} \begin{pmatrix} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \Rightarrow \mathbf{X}(\mathbf{V}_1, \mathbf{V}_2) = (\mathbf{U}_1, \mathbf{U}_2) \begin{pmatrix} \mathbf{D}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix},$$

desde donde sigue que

$$\mathbf{X} \mathbf{V}_1 = \mathbf{U}_1 \mathbf{D}_1, \quad \mathbf{X} \mathbf{V}_2 = \mathbf{0}.$$

Es decir, SVD permite “detectar” la dependencia lineal.

Una vez que hemos detectado que estamos en presencia de colinealidad, podemos sobrellevarla usando, por ejemplo, métodos de estimación sesgados. A continuación revisaremos dos procedimientos, *regresión por componentes principales* y el *estimador ridge*.

3.1.1. Regresión por componentes principales. Considere la descomposición espectral de $\mathbf{X}^\top \mathbf{X}$, dada por

$$\mathbf{X}^\top \mathbf{X} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top = (\mathbf{U}_1, \mathbf{U}_2) \begin{pmatrix} \mathbf{\Lambda}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_2 \end{pmatrix} \begin{pmatrix} \mathbf{U}_1^\top \\ \mathbf{U}_2^\top \end{pmatrix},$$

donde $\mathbf{\Lambda}_1 = \text{diag}(\lambda_1, \dots, \lambda_r)$ y $\mathbf{\Lambda}_2 = \text{diag}(\lambda_{r+1}, \dots, \lambda_p)$, mientras que $\mathbf{U} = (\mathbf{U}_1, \mathbf{U}_2)$ es matriz ortogonal. Esto lleva a la siguiente definición.

RESULTADO 3.2 (Estimador componentes principales). *Bajo los supuestos del modelo lineal en A1-A4*, el estimador componentes principales para β puede ser escrito como*

$$\hat{\beta}_r = U_1(U_1^\top X^\top XU_1)^{-1}U_1^\top X^\top Y = U_1\Lambda_1^{-1}U_1^\top X^\top Y.$$

DEMOSTRACIÓN. Por la ortogonalidad de $U = (U_1, U_2)$, sigue que

$$U_1^\top U_1 = I_r, \quad U_2^\top U_2 = I_{p-r}, \quad U_1 U_1^\top + U_2 U_2^\top = I_p,$$

y $U_1^\top U_2 = \mathbf{0}$. Ahora,

$$(X^\top X)^{-1} = U\Lambda^{-1}U^\top = U_1\Lambda_1^{-1}U_1^\top + U_2\Lambda_2^{-1}U_2^\top.$$

Usando que $U_1^\top U_2 = \mathbf{0}$ ($= U_2^\top U_1$), obtenemos

$$U_2^\top (X^\top X)^{-1} U_2 = U_2^\top (U_1\Lambda_1^{-1}U_1^\top + U_2\Lambda_2^{-1}U_2^\top) U_2 = \Lambda_2^{-1}$$

De este modo, $[U_2^\top (X^\top X)^{-1} U_2]^{-1} = \Lambda_2$, lo que permite escribir

$$\begin{aligned} (X^\top X)^{-1} U_2 [U_2^\top (X^\top X)^{-1} U_2]^{-1} U_2^\top (X^\top X)^{-1} \\ = (U_1\Lambda_1^{-1}U_1^\top + U_2\Lambda_2^{-1}U_2^\top) U_2 \Lambda_2 U_2^\top (U_1\Lambda_1^{-1}U_1^\top + U_2\Lambda_2^{-1}U_2^\top) \\ = U_2 \Lambda_2^{-1} U_2^\top. \end{aligned}$$

Es decir,

$$(X^\top X)^{-1} - (X^\top X)^{-1} U_2 [U_2^\top (X^\top X)^{-1} U_2]^{-1} U_2^\top (X^\top X)^{-1} = U_1 \Lambda_1^{-1} U_1^\top.$$

Como $U_1^\top (X^\top X)^{-1} U_1 = \Lambda_1$. Resulta

$$\begin{aligned} \hat{\beta}_r &= [(X^\top X)^{-1} - (X^\top X)^{-1} U_2 [U_2^\top (X^\top X)^{-1} U_2]^{-1} U_2^\top (X^\top X)^{-1}] X^\top Y \\ &= U_1 (U_1^\top X^\top XU_1)^{-1} U_1^\top X^\top Y, \end{aligned}$$

lo que concluye la prueba. \square

OBSERVACIÓN 3.3. El estimador PC es un caso particular del *estimador restringido* con respecto a:

$$U_2^\top \beta = \mathbf{0}.$$

En efecto, $\hat{\beta}_r$ depende del ‘parámetro’ r , basta notar que

$$\hat{\beta} = (X^\top X)^{-1} X^\top Y = (U_1\Lambda_1^{-1}U_1^\top + U_2\Lambda_2^{-1}U_2^\top) X^\top Y.$$

De este modo podemos interpretar $\hat{\beta}_r$ como una modificación del estimador OLS que *desconsidera* $U_2\Lambda_2^{-1}U_2^\top$.

Una alternativa para seleccionar r , es utilizar el test F . Suponga r fijo y considere $H_0 : U_2^\top \beta = \mathbf{0}$. Tenemos el estadístico

$$F = \left(\frac{n-p}{p-r} \right) \frac{(\hat{\beta} - \hat{\beta}_r)^\top X^\top X (\hat{\beta} - \hat{\beta}_r)}{Y^\top (I - H) Y}.$$

Si para un nivel α tenemos

$$F \geq F_{1-\alpha}(p-r, n-p).$$

Entonces, rechazamos H_0 y podemos seleccionar r un poco más pequeño.

OBSERVACIÓN 3.4. Sobre el estimador por componentes principales, debemos resaltar:

- No hay manera de verificar si las restricciones son satisfechas, de modo que este estimador es *sesgado*.
- Deseamos escoger r tan pequeño como posible para solucionar el problema de colinealidad y tan grande para no introducir mucho sesgo.

3.1.2. Estimador ridge. Hoerl y Kennard (1970) propusieron usar el *estimador ridge*, definido como:

$$\hat{\beta}_k = (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y}, \quad \lambda \geq 0$$

donde k es conocido como *parámetro ridge*. Evidentemente, podemos escribir el estimador ridge en términos del estimador de mínimos cuadrados, como:

$$\hat{\beta}_k = (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \hat{\beta}.$$

De este modo,

$$\mathbf{E}(\hat{\beta}_k) = (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \beta,$$

para $k \neq 0$, tenemos que $\hat{\beta}_k$ es un estimador sesgado. Mientras que su error cuadrático medio es dado por:

$$\text{MSE} = \mathbf{E}\{\|\hat{\beta}_k - \beta\|^2\} = \text{tr Cov}(\hat{\beta}_k) + \|\mathbf{E}(\hat{\beta}_k) - \beta\|^2.$$

En efecto,

$$\begin{aligned} \text{Cov}(\hat{\beta}_k) &= (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \text{Cov}(\hat{\beta}) \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1}. \end{aligned}$$

Además,

$$\begin{aligned} \text{bias}(\hat{\beta}_k, \hat{\beta}) &= \mathbf{E}(\hat{\beta}_k) - \beta = (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \beta - \beta \\ &= (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} [\mathbf{X}^\top \mathbf{X} \beta - (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})\beta] \\ &= -k(\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \beta. \end{aligned}$$

Considere la SVD de $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top$, de este modo $\mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top$, y podemos escribir

$$\begin{aligned} \text{Cov}(\hat{\beta}_k) &= \sigma^2 \mathbf{V} (\mathbf{D}^2 + k\mathbf{I})^{-1} \mathbf{V}^\top \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top \mathbf{V} (\mathbf{D}^2 + k\mathbf{I})^{-1} \mathbf{V}^\top \\ &= \sigma^2 \mathbf{V} (\mathbf{D}^2 + k\mathbf{I})^{-2} \mathbf{D}^2 \mathbf{V}^\top. \end{aligned}$$

De este modo,

$$\text{tr Cov}(\hat{\beta}_k) = \sigma^2 \text{tr}(\mathbf{D}^2 + k\mathbf{I})^{-2} \mathbf{D}^2 = \sigma^2 \sum_{i=1}^p \frac{\delta_i^2}{(\delta_i^2 + k)^2},$$

donde $\delta_1, \dots, \delta_p$ son los valores singulares de \mathbf{X} . Finalmente,

$$\text{MSE} = \sigma^2 \sum_{i=1}^p \frac{\delta_i^2}{(\delta_i^2 + k)^2} + k^2 \beta^\top (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-2} \beta.$$

El estimador ridge tiene varias interpretaciones interesantes, por ejemplo:

- (a) Es posible caracterizar $\hat{\beta}_k$ como solución del problema *regularizado*:

$$\min_{\beta} Q(\beta, k), \quad Q(\beta, k) = \|\mathbf{Y} - \mathbf{X}\beta\|^2 + k \|\beta\|^2, \quad (3.2)$$

que puede ser expresado de forma equivalente como

$$\min_{\beta} Q(\beta), \quad \text{sujeto a: } \|\beta\|^2 \leq r^2,$$

y en este contexto, k corresponde a un multiplicador de Lagrange.

- (b) Considere el modelo de regresión con *datos aumentados*:

$$\mathbf{Y}_a = \mathbf{X}_a \beta + \epsilon_a, \quad \epsilon_a \sim \mathbf{N}_{n+p}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

donde

$$\mathbf{Y}_a = \begin{pmatrix} \mathbf{Y} \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{X}_a = \begin{pmatrix} \mathbf{X} \\ \sqrt{k} \mathbf{I}_p \end{pmatrix}, \quad \epsilon_a = \begin{pmatrix} \epsilon \\ u \end{pmatrix}.$$

El interés recae en escoger algún $k > 0$ tal que la matriz de diseño \mathbf{X}_a tenga número condición $\kappa(\mathbf{X}_a)$ acotado.

OBSERVACIÓN 3.5. El tipo de regularización introducida en (3.2) es conocida como *regularización de Tikhonov*.¹

RESULTADO 3.6. *Suponga que los supuestos del modelo lineal en A1-A4*, son satisfechos. Entonces,*

$$\|\hat{\beta}_{k_2}\|^2 < \|\hat{\beta}_{k_1}\|^2,$$

siempre que $0 \leq k_1 < k_2$.

DEMOSTRACIÓN. Tenemos $\hat{\beta}_k = (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \hat{\beta}$. De este modo,

$$\|\hat{\beta}_k\|^2 = \hat{\beta}^\top \mathbf{M}_k \hat{\beta}, \quad \mathbf{M}_k = \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-2} \mathbf{X}^\top \mathbf{X}.$$

Basado en la SVD de \mathbf{X} , tenemos

$$\begin{aligned} \mathbf{M}_k &= \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top (\mathbf{V} \mathbf{D}^2 \mathbf{V}^\top + k \mathbf{V} \mathbf{V}^\top)^{-2} \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top \\ &= \mathbf{V} (\mathbf{D}^2 + k\mathbf{I})^{-2} \mathbf{D}^4 \mathbf{V}^\top = \mathbf{V} \mathbf{\Gamma}_k \mathbf{V}^\top, \end{aligned}$$

con

$$\mathbf{\Gamma}_k = \text{diag} \left(\frac{\delta_1^4}{(\delta_1^2 + k)^2}, \dots, \frac{\delta_p^4}{(\delta_p^2 + k)^2} \right).$$

De ahí que, si $0 \leq k_1 < k_2$, entonces

$$\mathbf{M}_{k_1} - \mathbf{M}_{k_2} \geq \mathbf{0},$$

lo que lleva a $\hat{\beta}^\top \mathbf{M}_{k_2} \hat{\beta} < \hat{\beta}^\top \mathbf{M}_{k_1} \hat{\beta}$, siempre que $\hat{\beta} \neq \mathbf{0}$. □

OBSERVACIÓN 3.7. Note que $\lim_{k \rightarrow \infty} \|\hat{\beta}_k\|^2 = 0$ y de ahí que

$$\lim_{k \rightarrow \infty} \hat{\beta}_k = \mathbf{0}. \quad (3.3)$$

Dado que $\hat{\beta}_k = \mathbf{W}_k \hat{\beta}$ con $\mathbf{W}_k = (\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X}$. La propiedad en (3.3) ha llevado a que el estimador ridge sea considerado como un *estimador shrinkage*, en cuyo caso

$$\mathbf{W}_k = (\mathbf{I}_p + k(\mathbf{X}^\top \mathbf{X})^{-1})^{-1}, \quad k \geq 0,$$

es llamada matrix ridge-shrinking.

¹Razón por la que k en ocasiones es llamado *parámetro de regularización*.

Se ha propuesto diversos estimadores de k , que buscan seleccionar un $\hat{\beta}_{k_{\text{opt}}}$ que reduzca su MSE. Algunas de estas alternativas son:

(a) La propuesta de [Hoerl, Kennard y Baldwin \(1975\)](#), dada por

$$\hat{k}_{\text{HKB}} = \frac{ps^2}{\|\hat{\beta}\|^2}.$$

(b) Estimador sugerido por [Lawless y Wang \(1976\)](#)

$$\hat{k}_{\text{LW}} = \frac{ps^2}{\hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta}}.$$

(c) Mientras que usando argumentos bayesianos, [Lindley y Smith \(1972\)](#) propusieron usar:

$$\hat{k}_{\text{LS}} = \frac{(n-p)(p+2)}{(n+2)} \frac{s^2}{\|\hat{\beta}\|^2}.$$

Por otro lado, [Golub, Heath y Wahba \(1979\)](#) han sugerido seleccionar el parámetro ridge usando *validación cruzada generalizada* (GCV), la que minimiza el criterio

$$V(k) = \frac{1}{n} \frac{\sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \hat{\beta}_k)^2}{\{1 - \text{tr}(\mathbf{H}(k))/n\}^2},$$

donde

$$\mathbf{H}(k) = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + k\mathbf{I})^{-1} \mathbf{X}^\top.$$

Es facil notar que $\hat{\mathbf{Y}}(k) = \mathbf{H}(k)\mathbf{Y}$. En este contexto se ha definido

$$\text{edf} = \text{tr} \mathbf{H}(k),$$

como el *número de parámetros efectivos*. En efecto, para $k = 0$, sigue que $\text{edf} = p$.

Considere la SVD de $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ y escriba el modelo en su *forma canónica*:

$$\mathbf{Z} = \mathbf{D}\boldsymbol{\alpha} + \mathbf{u}, \quad \mathbf{u} = \mathbf{U}^\top \boldsymbol{\epsilon} \sim \mathcal{N}_p(\mathbf{0}, \sigma^2 \mathbf{I}),$$

donde $\mathbf{Z} = \mathbf{U}^\top \mathbf{Y}$, $\boldsymbol{\alpha} = \mathbf{V}^\top \boldsymbol{\beta}$. De este modo, es fácil notar que

$$\hat{\boldsymbol{\alpha}}_k = (\mathbf{D}^2 + k\mathbf{I})^{-1} \mathbf{D}\mathbf{Z},$$

es decir,

$$\hat{\alpha}_{k,j} = \frac{\delta_j z_j}{\delta_j^2 + k}, \quad j = 1, \dots, p.$$

Por otro lado,

$$\begin{aligned} \text{edf} &= \text{tr} \mathbf{H}(k) = \text{tr} \mathbf{U}\mathbf{D}\mathbf{V}^\top (\mathbf{V}\mathbf{D}^2\mathbf{V}^\top + k\mathbf{I})^{-1} \mathbf{V}\mathbf{D}\mathbf{U}^\top \\ &= \text{tr} \mathbf{U}(\mathbf{D}^2 + k\mathbf{I})^{-1} \mathbf{D}^2 \mathbf{U}^\top = \text{tr}(\mathbf{D}^2 + k\mathbf{I})^{-1} \mathbf{D}^2 \\ &= \sum_{j=1}^p \frac{\delta_j^2}{\delta_j^2 + k}. \end{aligned}$$

Además,

$$\hat{\mathbf{Y}}(k) = \mathbf{X}\hat{\boldsymbol{\beta}}_k = \mathbf{U}\mathbf{D}\mathbf{V}^\top \hat{\boldsymbol{\beta}}_k = \mathbf{U}\mathbf{D}\hat{\boldsymbol{\alpha}}_k.$$

Lo que permite escribir

$$V(k) = \frac{1}{n} \frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_k\|^2}{\{\text{tr}(\mathbf{I} - \mathbf{H}(k))/n\}^2} = \frac{\|\mathbf{Y} - \mathbf{U}\mathbf{D}\hat{\boldsymbol{\alpha}}_k\|^2/n}{(1 - \text{edf}/n)^2}.$$

OBSERVACIÓN 3.8. Las consideraciones anteriores ofrecen un procedimiento sencillo para evaluar $V(k)$.

En términos del modelo canónico, también podemos escribir:

$$\hat{k}_{\text{HKB}} = \frac{ps^2}{\|\hat{\alpha}\|^2}, \quad \hat{k}_{\text{LW}} = \frac{ps^2}{\|\mathbf{Z}\|^2},$$

con $\hat{\alpha} = \mathbf{D}^{-1}\mathbf{Z}$ ($= \hat{\alpha}_0$), $\mathbf{Z} = \mathbf{U}^\top \mathbf{Y}$ y $s^2 = \|\mathbf{Y} - \mathbf{UD}\hat{\alpha}\|^2/(n-p)$.

EJEMPLO 3.9. [Woods, Steinour y Starke \(1932\)](#) consideraron datos desde un estudio experimental relacionando la emisión de calor durante la producción y endurecimiento de 13 muestras de *cementos Portland*. Este estudio se enfoca en los cuatro compuestos para los clinkers desde los que se produce el cemento. La respuesta (Y) es la *emisión de calor* después de 180 días de curado, medido en calorías por gramo de cemento. Los regresores son los porcentajes de los cuatro compuestos principales: *aluminato tricálcico* (X_1), *silicato tricálcico* (X_2), *ferrito aluminato tetracálcico* (X_3) y *silicato dicálcico* (X_4).

Siguiendo a [Woods, Steinour y Starke \(1932\)](#) consideramos un modelo lineal sin intercepto (modelo homogéneo), cuyo número condición escalado es $\kappa(\mathbf{X}) = 9.432$, esto es, \mathbf{X} es bien condicionada (mientras que para las variables centradas, obtenemos $\kappa(\tilde{\mathbf{X}}) = 37.106$).

Por otro lado, [Hald \(1952\)](#), [Gorman y Toman \(1966\)](#) y [Daniel y Wood \(1980\)](#) adoptan un modelo con intercepto (modelo no homogéneo). En cuyo caso $\kappa(\mathbf{X}) = 249.578$, sugiriendo la presencia de colinealidad. El aumento en el número condición se debe a que existe una relación lineal aproximada, pues

$$x_1 + x_2 + x_3 + x_4 \approx 100,$$

de modo que incluir el intercepto causa una colinealidad severa. Podemos usar la rutina para regresión ridge disponible en la biblioteca `fastmatrix` ([Osorio y Ogueda, 2021](#))

```
# carga biblioteca 'fastmatrix' y base de datos
> library(fastmatrix)
> load("portland.rda")

# ajuste usando regresión ridge
> ridge(y ~ x1 + x2 + x3 + x4, data = portland, lambda = 10,
+       method = "grid")

Call:
ridge(formula = y ~ x1 + x2 + x3 + x4, data = portland,
      lambda = 10, method = "grid")

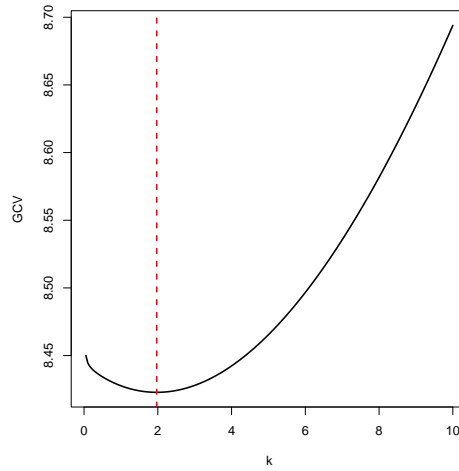
Coefficients:
(Intercept)          x1          x2          x3          x4
   0.08568    2.16549    1.15860    0.73845    0.48948

Optimal ridge parameter: 1.9598
Number of observations: 13
Effective number of parameters: 3.9796
Scale parameter estimate: 4.0553
```

La función `ridge` desde el paquete `fastmatrix`, permite también obtener los estimadores sugeridos por [Hoerl, Kennard y Baldwin \(1975\)](#) y [Lawless y Wang \(1976\)](#). Para el caso de los datos de cemento, tenemos:

$$\hat{k}_{\text{HKB}} = 0.0077, \quad \hat{k}_{\text{LW}} = 0.0032,$$

mientras que el valor óptimo de k obtenido mediante minimizar el criterio de validación cruzada generalizada, es $\hat{k}_{\text{opt}} = 1.9598$. En la siguiente figura, se presenta evaluación de la función $V(k)$ para una grilla de valores de k . Se ha indicado el mínimo \hat{k}_{opt} como una línea segmentada en color rojo,



Adicionalmente, es interesante llevar a cabo una comparativa usando el estimador mínimos cuadrados en el modelo sin intercepto. Considere la siguiente tabla de resumen de estimación:

Parámetro	Estimación OLS		Estimación ridge		
	homogéneo	no homogéneo	HKB	LW	GCV
β_0	—	62.4054	8.5870	17.1889	0.0855
β_1	2.1930	1.5511	2.1046	2.0162	2.1653
β_2	1.1533	0.5102	1.0648	0.9762	1.1586
β_3	0.7585	0.1019	0.6681	0.5776	0.7383
β_4	0.4863	-0.1441	0.3996	0.3127	0.4895
σ^2	4.0469	3.6818	4.0005	3.9478	5.0902
k	—	—	0.0077	0.0032	1.9716
edf	4.0000	5.0000	4.1369	4.2749	3.9795
κ	9.4325	249.5783	92.4131	130.8854	10.7852

Desde la tabla se aprecia que la elección del parámetro ridge usando el método de validación cruzada lleva a un número condición bastante pequeño, y en efecto, los resultados de estimación son muy cercanos al modelo homogéneo.

3.2. Errores correlacionados

El objetivo de esta sección es revisar el supuesto de *homogeneidad de varianzas*. De este modo, nos enfocaremos en el siguiente modelo:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

donde $\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$, y

$$\text{Cov}(\boldsymbol{\epsilon}) = \mathbf{E}(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top) = \sigma^2\boldsymbol{\Omega}, \quad \boldsymbol{\Omega} > \mathbf{0}.$$

Primeramente vamos a suponer que $\boldsymbol{\Omega}$ es conocida y sea $\boldsymbol{\Omega} = \mathbf{B}\mathbf{B}^\top$, con \mathbf{B} matriz no singular $n \times n$. Note que

$$\mathbf{B}^{-1}\mathbf{Y} = \mathbf{B}^{-1}\mathbf{X}\boldsymbol{\beta} + \mathbf{B}^{-1}\boldsymbol{\epsilon}.$$

haciendo $\mathbf{Y}_* = \mathbf{B}^{-1}\mathbf{Y}$, $\mathbf{X}_* = \mathbf{B}^{-1}\mathbf{X}$ y $\boldsymbol{\epsilon}_* = \mathbf{B}^{-1}\boldsymbol{\epsilon}$. Entonces $\mathbf{E}(\boldsymbol{\epsilon}_*) = \mathbf{0}$, y

$$\text{Cov}(\boldsymbol{\epsilon}_*) = \mathbf{B}^{-1} \text{Cov}(\boldsymbol{\epsilon}) \mathbf{B}^{-\top} = \sigma^2 \mathbf{B}^{-1} \mathbf{B} \mathbf{B}^\top \mathbf{B}^{-\top} = \sigma^2 \mathbf{I}.$$

Es decir, el modelo transformado

$$\mathbf{Y}_* = \mathbf{X}_*\boldsymbol{\beta} + \boldsymbol{\epsilon}_*,$$

satisface las condiciones A1-A4. Así,

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{GLS}} &= (\mathbf{X}_*^\top \mathbf{X}_*)^{-1} \mathbf{X}_*^\top \mathbf{Y}_* \\ &= (\mathbf{X}^\top \mathbf{B}^{-\top} \mathbf{B}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{B}^{-\top} \mathbf{B}^{-1} \mathbf{Y} \\ &= (\mathbf{X}^\top (\mathbf{B}\mathbf{B}^\top)^{-1} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{B}\mathbf{B}^\top)^{-1} \mathbf{Y} \\ &= (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{Y}. \end{aligned}$$

Es fácil mostrar que

$$\begin{aligned} \mathbf{E}(\hat{\boldsymbol{\beta}}_{\text{GLS}}) &= (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{E}(\mathbf{Y}) = \boldsymbol{\beta} \\ \text{Cov}(\hat{\boldsymbol{\beta}}_{\text{GLS}}) &= \sigma^2 (\mathbf{X}_*^\top \mathbf{X}_*)^{-1} = \sigma^2 (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1}. \end{aligned}$$

Además, podemos definir el vector de residuos,

$$\mathbf{e}_* = \mathbf{Y}_* - \mathbf{X}_* \hat{\boldsymbol{\beta}}_{\text{GLS}} = \mathbf{B}^{-1} \mathbf{Y} - \mathbf{B}^{-1} \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}} = \mathbf{B}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}}),$$

lo que permite escribir la suma de cuadrados de residuos, como:

$$\begin{aligned} Q_\Omega(\hat{\boldsymbol{\beta}}_{\text{GLS}}) &= \|\mathbf{e}_*\|^2 = (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}})^\top \mathbf{B}^{-\top} \mathbf{B}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}}) \\ &= (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}})^\top \boldsymbol{\Omega}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}}). \end{aligned}$$

Adicionalmente, podemos escribir:

$$\begin{aligned} \mathbf{e}_* &= \mathbf{B}^{-1} \mathbf{Y} - \mathbf{B}^{-1} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{Y} \\ &= \mathbf{B}^{-1} \mathbf{Y} - \mathbf{B}^{-1} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{B}^{-\top} \mathbf{B}^{-1} \mathbf{Y}, \end{aligned}$$

de ahí que

$$\mathbf{e}_* = (\mathbf{I} - \mathbf{H}_\Omega) \mathbf{Y}_*,$$

con

$$\mathbf{H}_\Omega = \mathbf{B}^{-1} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{B}^{-\top}.$$

OBSERVACIÓN 3.10. Evidentemente $\mathbf{H}_\Omega^\top = \mathbf{H}_\Omega$ y $\mathbf{H}_\Omega^2 = \mathbf{H}_\Omega$. De este modo la suma de cuadrados residual adopta la forma:

$$Q_\Omega(\hat{\boldsymbol{\beta}}_{\text{GLS}}) = \|\mathbf{e}_*\|^2 = \mathbf{Y}_*^\top (\mathbf{I} - \mathbf{H}_\Omega) \mathbf{Y}_*.$$

Desafortunadamente, en general la matriz $\mathbf{\Omega}$ no es conocida y requiere ser estimada. Si $\hat{\mathbf{\Omega}}$ es un estimador de $\mathbf{\Omega}$, entonces

$$\hat{\beta}_{\text{EGLS}} = (\mathbf{X}^\top \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{\Omega}}^{-1} \mathbf{Y}.$$

Se debe notar que las propiedades de $\hat{\beta}_{\text{EGLS}}$ son difíciles de caracterizar.

Un caso particular importante, corresponde a *mínimos cuadrados ponderados (WLS)* en cuyo caso $\text{Cov}(\epsilon) = \sigma^2 \mathbf{W}^{-1}$ donde $\mathbf{W} = \text{diag}(\omega_1, \dots, \omega_n)$ con $\omega_i > 0, \forall i$. De este modo,

$$\hat{\beta}_{\text{WLS}} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}.$$

Note que este estimador es solución del problema

$$\min_{\beta} Q_W(\beta), \quad Q_W(\beta) = (\mathbf{Y} - \mathbf{X}\beta)^\top \mathbf{W} (\mathbf{Y} - \mathbf{X}\beta).$$

Bajo el supuesto $\epsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{W}^{-1})$, con $\mathbf{W} = \text{diag}(\omega_1, \dots, \omega_n)$. Es fácil notar que el estimador $\hat{\beta}_{\text{WLS}}$ minimiza la función

$$Q_W(\beta) = \sum_{i=1}^n \omega_i (Y_i - \mathbf{x}_i^\top \beta)^2.$$

Además, el estimador $\hat{\beta}_{\text{WLS}}$ resuelve las siguientes ecuaciones de estimación

$$\mathbf{X}^\top \mathbf{W} (\mathbf{Y} - \mathbf{X}\beta) = \mathbf{0},$$

o equivalentemente

$$\sum_{i=1}^n \omega_i (Y_i - \mathbf{x}_i^\top \beta) \mathbf{x}_i = \mathbf{0}.$$

Mientras que el estimador ML para σ^2 asume la forma:

$$\hat{\sigma}_{\text{WLS}}^2 = \frac{1}{n} \sum_{i=1}^n \omega_i (Y_i - \mathbf{x}_i^\top \hat{\beta}_{\text{WLS}})^2.$$

OBSERVACIÓN 3.11. El problema anterior puede ser resuelto usando OLS en el siguiente ‘problema modificado’

$$\mathbf{W}^{1/2} \mathbf{Y}, \quad \mathbf{W}^{1/2} \mathbf{X}.$$

Consideraciones computacionales sobre GLS y WLS son dadas por ejemplo en el Capítulo 4 de [Björck \(1996\)](#).

3.2.1. Estimación de funciones de varianza. El objetivo de esta sección es considerar modelos de regresión heterocedásticos, tal que

$$\mathbf{E}(Y_i) = \mu_i, \quad \text{var}(Y_i) = \sigma^2 g^2(\mathbf{z}_i; \mu_i, \phi), \quad i = 1, \dots, n,$$

donde $\mu_i = \mathbf{x}_i^\top \beta$, \mathbf{x}_i y \mathbf{z}_i representan vectores de covariables (que podrían ser iguales), β son coeficientes de regresión, g es función de varianza (que permite modelar la heterogeneidad), $\sigma^2 > 0$ y ϕ son parámetros de escala desconocidos.

EJEMPLO 3.12. Considere

$$\text{var}(Y_i) = \sigma^2 \{g(\mathbf{x}_i; \beta)\}^{2\phi}, \quad \phi > 0.$$

Si suponemos $g(\mathbf{x}_i; \beta) = \mathbf{x}_i^\top \beta$ y $\phi = \frac{1}{2}$ tenemos la *estructura de varianza Poisson*, mientras que $\phi = 1$ es de *tipo-gama*.

Ejemplos habituales de funciones de varianza son los siguientes:

(a) Función de varianza cuadrática

$$\sigma g(\mathbf{z}_i; \mu_i, \phi) = 1 + \phi_1 z_{1i} + \phi_2 z_{2i}^2.$$

(b) Modelo potencia extendido

$$\text{var}(Y_i) = \sigma^2 (\phi_1 + \phi_2 \mu_i^{\phi_3}).$$

(c) Modelo exponencial

$$\text{var}(Y_i) = \sigma^2 \exp(2\phi \mu_i).$$

(d) También puede depender de ϕ según un predictor lineal

$$\text{var}(Y_i) = \sigma^2 \exp(2\mathbf{z}_i^\top \boldsymbol{\phi}).$$

Primeramente, supondremos el modelo

$$\mathbb{E}(Y_i) = \mu_i = \mathbf{x}_i^\top \boldsymbol{\beta}, \quad \text{var}(Y_i) = \sigma^2 g^2(\mathbf{z}_i; \mu_i, \phi),$$

con ϕ conocido. El método WLS sugiere un procedimiento natural para la estimación de varianzas heterogéneas usando una estrategia de *mínimos cuadrados iterativamente ponderados (IWLS)*.

Algoritmo 4: IWLS para estimación de varianzas

```

1 begin
2   Considerar una estimación inicial para  $\boldsymbol{\beta}$ , digamos  $\boldsymbol{\beta}^{(0)}$ , resolviendo
      
$$\sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}) \mathbf{x}_i = \mathbf{0}.$$

      y hacemos  $k \leftarrow 0$ .
3   Construir pesos
      
$$\omega_i^{(k)} = 1/g^2(\mathbf{z}_i; \mu_i^{(k)}, \phi), \quad \mu_i^{(k)} = \mathbf{x}_i^\top \boldsymbol{\beta}^{(k)}.$$

4   Actualizar  $\boldsymbol{\beta}^{(k+1)}$ , resolviendo
      
$$\sum_{i=1}^n \omega_i^{(k)} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}) \mathbf{x}_i = \mathbf{0},$$

      hacer  $k \leftarrow k + 1$  y volver a Paso 3.
5 end
```

Evidentemente el Paso 4 del Algoritmo 4, debe ser resuelta usando mínimos cuadrados ponderados.

En el algoritmo anterior σ^2 puede ser estimado por analogía a WLS. Específicamente, podemos considerar

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\omega}_i (Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2,$$

donde $\hat{\omega}_i$ ($i = 1, \dots, n$) y $\hat{\boldsymbol{\beta}}$ representan los valores de ω_i y $\boldsymbol{\beta}$ a la convergencia del Algoritmo 4.

Alternativamente, es posible incorporar una etapa adicional al **Paso 4** del Algoritmo 4 como:

$$\sigma^{2(k+1)} = \frac{1}{n} \sum_{i=1}^n \omega_i^{(k)} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}^{(k+1)})^2.$$

Claramente, no siempre es posible especificar valores para $\boldsymbol{\phi}$. Para simplificar la exposición considere

$$Y_i \stackrel{\text{ind}}{\sim} \mathbf{N}(\mu_i, \sigma^2 g^2(\mathbf{z}_i; \mu_i, \boldsymbol{\phi})), \quad i = 1, \dots, n,$$

con $\mu_i = \mathbf{x}_i^\top \boldsymbol{\beta}$. Defina la matriz diagonal

$$\mathbf{G} = \mathbf{G}(\boldsymbol{\beta}, \boldsymbol{\phi}) = \text{diag}(g^2(\mathbf{z}_1; \mu_1, \boldsymbol{\phi}), \dots, g^2(\mathbf{z}_n; \mu_n, \boldsymbol{\phi})).$$

De este modo, el modelo anterior puede ser escrito como

$$\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{G}),$$

con función de densidad conjunta

$$f(\mathbf{y}) = |2\pi\sigma^2 \mathbf{G}|^{-1/2} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{G}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\}.$$

La parte relevante de la función de log-verosimilitud es

$$\ell_n(\boldsymbol{\theta}) = -\frac{1}{2} \log |\sigma^2 \mathbf{G}| - \frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{G}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}),$$

con $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \sigma^2, \boldsymbol{\phi}^\top)^\top$. Por la estructura diagonal de \mathbf{G} tenemos que

$$\begin{aligned} \mathbf{G}^{-1} &= \text{diag}(g^{-2}(\mathbf{z}_1; \mu_1, \boldsymbol{\phi}), \dots, g^{-2}(\mathbf{z}_n; \mu_n, \boldsymbol{\phi})), \\ \log |\sigma^2 \mathbf{G}| &= \sum_{i=1}^n \log \sigma^2 g^2(\mathbf{z}_i; \mu_i, \boldsymbol{\phi}). \end{aligned}$$

Esto permite escribir la función de log-verosimilitud como:

$$\begin{aligned} \ell_n(\boldsymbol{\theta}) &= -\frac{1}{2} \sum_{i=1}^n \log \sigma^2 g^2(\mathbf{z}_i; \mu_i, \boldsymbol{\phi}) - \frac{1}{2\sigma^2} \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{g^2(\mathbf{z}_i; \mu_i, \boldsymbol{\phi})} \\ &= -\frac{1}{2} \sum_{i=1}^n \left\{ \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{\sigma^2 g^2(\mathbf{z}_i; \mu_i, \boldsymbol{\phi})} + \log \sigma^2 g^2(\mathbf{z}_i; \mu_i, \boldsymbol{\phi}) \right\} \end{aligned}$$

La estimación de parámetros se puede desarrollar alternando las siguientes dos etapas:

1. Para una estimación preliminar $\boldsymbol{\beta}^{(k)}$ de $\boldsymbol{\beta}$ minimizar con relación a $\boldsymbol{\phi}$ y σ^2 , la función de *log-verosimilitud perfilada*:

$$\ell_*(\boldsymbol{\beta}^{(k)}, \sigma^2, \boldsymbol{\beta}) = -\frac{1}{2} \sum_{i=1}^n \left\{ \frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}^{(k)})^2}{\sigma^2 g^2(\mathbf{z}_i; \mu_i^{(k)}, \boldsymbol{\phi})} + \log \sigma^2 g^2(\mathbf{z}_i; \mu_i^{(k)}, \boldsymbol{\phi}) \right\},$$

con $\mu_i^{(k)} = \mathbf{x}_i^\top \boldsymbol{\beta}^{(k)}$. Diferenciando con relación a $\boldsymbol{\phi}$ y σ^2 lleva a las ecuaciones de estimación:

$$\sum_{i=1}^n \frac{1}{g^4(\mathbf{z}_i; \mu_i^{(k)}, \boldsymbol{\phi})} \{r_i^2 - \sigma^2 g^2(\mathbf{z}_i; \mu_i^{(k)}, \boldsymbol{\phi})\} \mathbf{q}(\mu_i^{(k)}, \sigma, \boldsymbol{\phi}) = \mathbf{0}, \quad (3.4)$$

donde

$$\mathbf{q}(\mu_i^{(k)}, \sigma, \boldsymbol{\phi}) = g^2(\mu_i^{(k)}, \boldsymbol{\phi}) \begin{pmatrix} 1/\sigma \\ \mathbf{u}(\mu_i^{(k)}, \boldsymbol{\phi}) \end{pmatrix},$$

y $\mathbf{u}(\mu_i^{(k)}, \phi)$ representa el vector de derivadas de $\log g(\mathbf{z}_i, \mu_i, \phi)$ con relación a ϕ , mientras que $r_i = Y_i - \mathbf{x}_i^\top \beta^{(k)}$.

2. Actualizar $\beta^{(k+1)}$ como la solución del problema mínimos cuadrados ponderados

$$\mathbf{X}^\top \mathbf{G}^{-1} (\mathbf{Y} - \mathbf{X}\beta) = \mathbf{0},$$

donde $\mathbf{G} = \mathbf{G}(\beta^{(k)}, \phi^{(k+1)})$.

Inspección de la Ecuación (3.4) permite notar que la estimación de σ y ϕ corresponde a WLS con *respuesta* r_i^2 , *función de regresión* $\sigma^2 g^2(\mathbf{z}_i; \mu_i^{(k)}, \phi)$, *parámetros de regresión* $(\sigma, \phi^\top)^\top$, *pesos* $g^{-4}(\mathbf{z}_i, \mu_i^{(k)}, \phi)$ y *gradiente* (matriz de diseño cuyas filas están dadas por) $\mathbf{q}(\mu_i^{(k)}, \sigma, \phi)$.

Detalles sobre las propiedades estadísticas del estimador obtenido por este procedimiento pueden ser hallados en [Davidian y Carroll \(1987\)](#) (ver también [Carroll y Ruppert, 1988](#)).

3.3. Transformaciones estabilizadoras de varianza

Para introducir ideas, considere

$$\mathbb{E}(Y) = \mu, \quad \text{var}(Y) = \sigma^2 h(Y),$$

y suponga la transformación $z = g(y)$ tal que la varianza de Z es aproximadamente independiente de μ . Adicionalmente, considere una expansión de primer orden de $g(y)$ en torno de μ , esto es

$$g(y) \approx g(\mu) + g'(\mu)(y - \mu).$$

De este modo,

$$\mathbb{E}(Z) \approx g(\mu)$$

$$\text{var}(Z) \approx \{g'(\mu)\}^2 \text{var}(Y - \mu) = \{g'(\mu)\}^2 \sigma^2 h(y)$$

Así, para determinar una transformación tal que $\text{var}(Z) = \sigma^2$ necesitamos que

$$g'(\mu) = \frac{1}{\sqrt{h(y)}},$$

o de forma análoga,

$$g(\mu) = \int \frac{d\mu}{\sqrt{h(\mu)}}.$$

En particular, se podría considerar la clase de transformaciones en que $h(y)$ es una potencia de μ . Algunos ejemplos son los siguientes:

$h(\mu)$	z	Descripción
μ^4	$1/y$	recíproco
μ^2	$\log y$	logarítmico
μ	\sqrt{y}	raíz cuadrada
$\mu(1 - \mu)$	$\sin^{-1}(\sqrt{y})$	seno inverso
$(1 - \mu^2)^2$	$\log(\frac{1+y}{1-y})$	correlación

[Box y Cox \(1964\)](#) sugirieron llevar a cabo la estimación ML para una clase general de transformaciones. El supuesto fundamental es que $Y > 0$ y que existe alguna

potencia de Y tal que su varianza es aproximadamente constante y satisface el supuesto de normalidad. En concreto, ellos consideraron la familia de transformaciones

$$Y(\lambda) = \begin{cases} \frac{Y^\lambda - 1}{\lambda}, & \lambda \neq 0, \\ \log(Y), & \lambda = 0, \end{cases}$$

el objetivo es obtener una estimación de λ desde los datos observados.

Sea $U_i = Y_i(\lambda)$ para $i = 1, \dots, n$, y asuma que $\mathbf{U} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ para alguna elección de λ . El Jacobiano de la transformación es dado por

$$J = \prod_{i=1}^n \frac{\partial Y_i(\lambda)}{\partial Y_i} = \prod_{i=1}^n Y_i^{\lambda-1} = \left(\prod_{i=1}^n Y_i^{1/n} \right)^{n(\lambda-1)} = G^{n(\lambda-1)},$$

donde G es la media geométrica de las observaciones. De este modo, la log-verosimilitud requerida adopta la forma:

$$\ell_n(\boldsymbol{\theta}, \lambda) = -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} Q_\lambda(\boldsymbol{\beta}) + \log J,$$

donde

$$Q_\lambda(\boldsymbol{\beta}) = \|\mathbf{Y}(\lambda) - \mathbf{X}\boldsymbol{\beta}\|^2.$$

Para λ fijado, tenemos

$$\hat{\boldsymbol{\beta}}(\lambda) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}(\lambda), \quad \hat{\sigma}^2(\lambda) = \frac{1}{n} \|\mathbf{Y}(\lambda) - \mathbf{X}\hat{\boldsymbol{\beta}}(\lambda)\|^2.$$

Esto lleva a la función de log-verosimilitud perfilada, dada por

$$\ell_*(\lambda) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log(\text{RSS}(\lambda)/n) - \frac{n}{2} + \log J,$$

donde

$$\text{RSS}(\lambda) = \mathbf{Y}^\top(\lambda)(\mathbf{I} - \mathbf{H})\mathbf{Y}(\lambda).$$

Considere por conveniencia,

$$\mathbf{Z}(\lambda) = \frac{\mathbf{Y}(\lambda)}{G^{\lambda-1}},$$

de este modo podemos escribir

$$\ell_*(\lambda) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log(\text{RSS}_Z(\lambda)/n) - \frac{n}{2},$$

con

$$\text{RSS}_Z(\lambda) = \mathbf{Z}^\top(\lambda)(\mathbf{I} - \mathbf{H})\mathbf{Z}(\lambda),$$

luego la maximización de $\ell_*(\lambda)$ es equivalente a la minimización de $\text{RSS}_Z(\lambda)$.

Típicamente se realiza la transformación $Z(\lambda)$ para un rango de valores de λ , se realiza el ajuste del modelo lineal y luego se examina aquél valor de λ que corresponde al valor más pequeño de $\text{RSS}_Z(\lambda)$ ²

Para los datos transformados la función Box-Cox asume la forma:

$$Z(\lambda) = \begin{cases} \frac{Y^\lambda - 1}{\lambda G^{\lambda-1}}, & \lambda \neq 0, \\ G \log(Y), & \lambda = 0, \end{cases}$$

obteniendo el valor $\hat{\lambda}$ se lleva a cabo el ajuste $\mathbf{E}(\mathbf{Z}) = \mathbf{X}\boldsymbol{\beta}$ con $\mathbf{Z} = \mathbf{Z}(\hat{\lambda})$.

²Usualmente basta considerar $-2 \leq \lambda \leq 2$ con incrementos no muy pequeños.

OBSERVACIÓN 3.13. Note que $\lambda = 1$ implica que el modelo no debe ser transformado.

Para llevar a cabo la estimación de parámetros usando el procedimiento de Box-Cox, podemos considerar el siguiente fragmento de código en R:

```
boxcox.lm <- function(x, y, lambda) {
  boxcox <- function(y, lambda) {
    n <- length(y)
    lambda <- rep(lambda, n)
    z <- ifelse(lambda != 0., (y^lambda - 1.) / lambda, log(y))
    z
  }
  n <- nrow(x)
  p <- ncol(x)
  k <- length(lambda)
  RSS <- rep(0, k)
  logLik <- rep(0, k)
  for (i in 1:k) {
    geom <- geomean(y)
    z <- boxcox(y, lambda = lambda[i])
    z <- z / geom^(lambda[i] - 1.)
    fm <- ols.fit(x, z, method = "sweep")
    RSS[i] <- fm$RSS
    logLik[i] <- -.5 * n * log(2 * pi)
      - .5 * n * log(RSS[i] / n) - .5 * n
  }
  idx <- order(RSS)[1]
  opt <- lambda[idx]
  obj <- list(lambda = lambda, RSS = RSS, logLik = logLik,
    opt = opt)
  obj
}
```

EJEMPLO 3.14. Para los datos de Forbes (ver [Weisberg, 2005](#)), llevamos a cabo selección de $\lambda \in [-2, 2]$ usando la transformación Box-Cox. Considere,

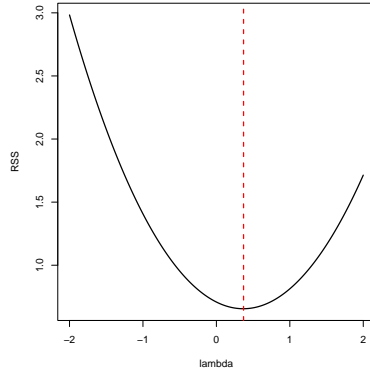
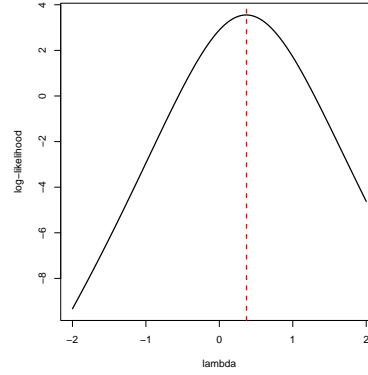
```
# carga datos desde biblioteca MASS
> library(MASS)
> data(forbes)

# ajuste preliminar
> library(fastmatrix)
> fm <- ols(pres ~ bp, data = forbes, x = TRUE, y = TRUE)

# extrae vector de respuestas y matriz de diseno
> x <- fm$x
> y <- fm$y

# interpreta script, crea 'grilla' e invoca método de ajuste
> source("boxcox.lm.R")
> lambda <- seq(-2, 2, by = 0.01)
> z <- boxcox.lm(x, y, lambda)
```

El valor óptimo basado en una grilla de valores para $\lambda \in \{-2.00, -1.99, \dots, 1.99, 2.00\}$, se obtuvo $\lambda_{\text{opt}} = 0.37$. Los siguientes gráficos presentan las funciones $\text{RSS}_Z(\lambda)$ y la log-verosimilitud perfilada:

(a) $\text{RSS}_Z(\lambda)$ (b) $\ell_*(\lambda)$

Adicionalmente, podemos realizar el análisis considerando algunos valores para λ así como eliminado la observación 12, en cuyo caso se obtuvo $\hat{\lambda}_{\text{opt}} = 0.10$. Considere el siguiente fragmento de código en R y el resumen de estimación:

```
# gráficos de RSS y log-likelihood perfilada
> plot(lambda, z$RSS, type = "l", ylab = "RSS", lwd = 2)
> abline(v = z$opt, col = "red", lwd = 2, lty = 2)
> plot(lambda, z$logLik, type = "l", ylab = "log-likelihood")
> abline(v = z$opt, col = "red", lwd = 2, lty = 2)

# removiendo dato 12
> y12 <- y[-12]
> x12 <- x[-12,]
> z <- boxcox.lm(x12, y12, lambda)

# ajustando diversos modelos
f0 <- ols(pres ~ bp, data = forbes)
f1 <- ols(log(pres) ~ bp, data = forbes)
f2 <- ols((pres^.37 - 1) / .37 ~ bp, data = forbes)
f3 <- ols((pres^.10 - 1) / .10 ~ bp, data = forbes)
```

Parámetro	λ			
	—	0.00	0.10	0.37
β_0	-81.0637	-0.9709	-1.9885	-7.6462
β_1	0.5229	0.0206	0.0285	0.0681
σ^2	0.0542	0.0001	0.0001	0.0008
RSS	0.8131	0.0011	0.0021	0.0114
$\ell(\hat{\theta})$	1.7186	57.5378	52.3791	37.9802

Es decir, el mejor modelo corresponde a la transformación Box-Cox con $\lambda = 0$ (esto es la transformación logarítmica). Adicionalmente es instructivo llevar a cabo el test de razón de verosimilitudes $H_0 : \lambda = 1$.

3.4. Análisis de residuos y leverages

Suponga el modelo lineal,

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

con los Supuestos A1-A4. El vector de residuos es dado por:

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y},$$

con $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$, es decir $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$ ($= \hat{\mathbf{E}}(\mathbf{Y})$). Bajo el supuesto de normalidad $\mathbf{Y} \sim \mathbf{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$, tenemos

$$\mathbf{e} \sim \mathbf{N}_n(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H})),$$

es decir

$$\mathbf{E}(e_i) = 0, \quad \text{var}(e_i) = \sigma^2(1 - h_{ii}), \quad \text{Cov}(e_i, e_j) = -\sigma^2 h_{ij}.$$

De ahí que los residuos tienen *varianzas diferentes y son correlacionados*. A continuación se introducirá definiciones de residuos estandarizados. Primeramente, considere que σ^2 es conocido, de este modo

$$z_i = \frac{e_i}{\sigma} \sim \mathbf{N}(0, 1).$$

De este modo, podemos definir el *residuo estandarizado* como:

$$r_i = \frac{e_i}{s\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n.$$

Cook y Weisberg (1982) mostraron que

$$\frac{r_i^2}{n - p} \sim \text{Beta}\left(\frac{1}{2}, \frac{n - p - 1}{2}\right).$$

Esto nos permite notar que,

$$\mathbf{E}(r_i) = 0, \quad \text{var}(r_i) = 1, \quad \text{Cov}(r_i, r_j) = \frac{-h_{ij}}{\sqrt{(1 - h_{ii})(1 - h_{jj})}}.$$

Por otro lado, considere el *residuo studentizado*:

$$t_i = \frac{e_i}{s_{(i)}\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n.$$

donde

$$s_{(i)}^2 = \frac{1}{n - p - 1} \sum_{j \neq i}^n (y_j - \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}_{(i)})^2,$$

$$\hat{\boldsymbol{\beta}}_{(i)} = (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^\top \mathbf{Y}_{(i)},$$

denotan los estimadores de σ^2 y $\boldsymbol{\beta}$ una vez que la i -ésima observación ha sido eliminada. Es decir basados en las siguientes matrices de datos,

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_{(i)} \\ \mathbf{x}_i^\top \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} \mathbf{Y}_{(i)} \\ Y_i \end{pmatrix}$$

Una interpretación interesante de t_i es que corresponde al estadístico t para probar la hipótesis $H_0 : \gamma = 0$ en el *modelo de salto en la media*, dado por:

$$Y_j = \mathbf{x}_j^\top \boldsymbol{\beta} + d_j \gamma + \epsilon_j, \quad j = 1, \dots, n,$$

donde $d_j = 1$ si $j = i$ y 0 en caso contrario.

El modelo puede ser escrito como:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{d}_i\gamma + \boldsymbol{\epsilon},$$

con $\mathbf{d}_i = (\mathbf{0}, 1, \mathbf{0})^\top$ con un cero en la i -ésima posición.

Lo anterior permite notar que $t_i \sim t(n - p - 1)$, y de este modo,

$$\mathbb{E}(t_i) = 0, \quad \text{var}(t_i) = \frac{n - p - 1}{n - p - 3} \approx 1.$$

Es decir, los residuos estandarizados y estudentizados tienen propiedades similares a los errores $\{\epsilon_1, \dots, \epsilon_n\}$.

Basado en la propiedad $\text{Cov}(\mathbf{e}, \hat{\mathbf{Y}}) = \sigma^2(\mathbf{I} - \mathbf{H})\mathbf{H} = \mathbf{0}$ se ha sugerido el diagrama de residuos versus valores predichos. Este tipo de herramientas gráficas permite verificar desvios evidentes del modelo. Adicionalmente, se ha propuesto la construcción de gráficos cuantil-cuantil (QQ-plot) con *envelopes* para evaluar el supuesto de normalidad (consultar [Atkinson, 1985](#), para más detalles).

Por otro lado, sabemos que

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}, \quad \mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top, \quad (3.5)$$

de ahí sigue que,

$$\hat{Y}_i = \sum_{j=1}^n h_{ij} Y_j = h_{ii} Y_i + \sum_{j \neq i} h_{ij} Y_j, \quad (3.6)$$

es decir el valor predicho es una combinación lineal de las respuestas observadas con pesos dados por los elementos de la matriz de proyección \mathbf{H} . A continuación llamaremos a los elementos diagonales h_{ii} , $i = 1, \dots, n$, como *leverages*. A continuación revisamos algunas de las propiedades fundamentales de la matriz \mathbf{H} :

PROPIEDAD 3.15. \mathbf{H} es simétrica e idempotente con $\text{rg}(\mathbf{H}) = \text{tr}(\mathbf{H}) = p$.

PROPIEDAD 3.16. Los elementos diagonales de \mathbf{H} están acotados, en efecto:

$$0 \leq h_{ii} \leq 1, \quad i = 1, \dots, n.$$

DEMOSTRACIÓN. Sabemos que

$$\hat{\mathbf{Y}} \sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{H}), \quad \mathbf{e} \sim \mathbf{N}(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H})).$$

De este modo,

$$\text{var}(\hat{Y}_i) = \sigma^2 h_{ii}, \quad \text{var}(e_i) = \sigma^2(1 - h_{ii}),$$

de ahí sigue el resultado.

Para otra demostración, ver Resultado [A.6](#) desde el Apéndice [A](#). □

PROPIEDAD 3.17. Tenemos,

$$h_{ii} = \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i, \quad i = 1, \dots, n.$$

Además,

$$\bar{h} = \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{1}{n} \text{tr}(\mathbf{H}) = \frac{p}{n}.$$

PROPIEDAD 3.18. Sea $\widetilde{\mathbf{X}}$ la matriz de datos centrados. En este caso, los elementos diagonales de $\widetilde{\mathbf{H}}$ están dados por

$$\widetilde{h}_{ii} = (\mathbf{x}_i - \bar{\mathbf{x}})^\top (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}), \quad i = 1, \dots, n.$$

Luego \widetilde{h}_{ii} es la distancia ponderada desde \mathbf{x}_i al *centroide* $\bar{\mathbf{x}}$.

Hoaglin y Welsch (1978) sugirieron que aquellas observaciones que exceden dos veces su promedio

$$h_{ii} > 2p/n \quad (= 2\bar{h})$$

indican un alto leverage. Mientras que Huber (1981) sugirió identificar observaciones tal que

$$h_{ii} > 0.5,$$

independiente de n o p . En la práctica se debe prestar atención a casos inusualmente grandes *con relación al resto* de h_{ii} 's.

PROPIEDAD 3.19. Desde Ecuación (3.5), sigue que

$$\frac{\partial \widehat{\mathbf{Y}}}{\partial \mathbf{Y}^\top} = \mathbf{H},$$

y en particular $\partial \widehat{Y}_i / \partial Y_i = h_{ii}$, para $i = 1, \dots, n$.

PROPIEDAD 3.20. Si el modelo *tiene intercepto*, entonces $\mathbf{H}\mathbf{1} = \mathbf{1}$.

DEMOSTRACIÓN. Considere $\mathbf{X} = (\mathbf{1}, \mathbf{X}_1)$. Sabemos que $\mathbf{H}\mathbf{X} = \mathbf{X}$, y de ahí que

$$\mathbf{H}(\mathbf{1}, \mathbf{X}_1) = (\mathbf{1}, \mathbf{X}_1),$$

y el resultado sigue. \square

3.5. Diagnóstico de por eliminación de casos

Suponga que se desea evaluar el efecto de eliminar una observación sobre la estimación de $\boldsymbol{\beta}$. En este caso, podemos considerar el *modelo de datos eliminados*

$$\mathbf{Y}_{(i)} = \mathbf{X}_{(i)}\boldsymbol{\beta} + \boldsymbol{\epsilon}_{(i)}, \quad (3.7)$$

de ahí que

$$\widehat{\boldsymbol{\beta}}_{(i)} = (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^\top \mathbf{Y}_{(i)}.$$

Un aspecto interesante de la estimación de parámetros en el modelo definido en la Ecuación (3.7) es que no requiere re-ajustar n modelos (uno por cada observación que hemos removido), sino que podemos escribir el estimador $\widehat{\boldsymbol{\beta}}_{(i)}$ en términos de información calculada para el modelo con *datos completos*. Para notar esto, considere:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_{(i)} \\ \mathbf{x}_i^\top \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} \mathbf{Y}_{(i)} \\ Y_i \end{pmatrix},$$

de ahí que

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} &= (\mathbf{X}_{(i)}^\top, \mathbf{x}_i) \begin{pmatrix} \mathbf{X}_{(i)} \\ \mathbf{x}_i^\top \end{pmatrix} = \mathbf{X}_{(i)}^\top \mathbf{X}_{(i)} + \mathbf{x}_i \mathbf{x}_i^\top, \\ \mathbf{X}^\top \mathbf{Y} &= (\mathbf{X}_{(i)}^\top, \mathbf{x}_i) \begin{pmatrix} \mathbf{Y}_{(i)} \\ Y_i \end{pmatrix} = \mathbf{X}_{(i)}^\top \mathbf{Y}_{(i)} + \mathbf{x}_i Y_i. \end{aligned}$$

Reagrupando, podemos escribir

$$\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)} = \mathbf{X}^\top \mathbf{X} (\mathbf{I} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top),$$

cuya matriz inversa es dada por

$$\begin{aligned} (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} &= (\mathbf{I} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top)^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \left\{ \mathbf{I} + \frac{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top}{1 - \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i} \right\} (\mathbf{X}^\top \mathbf{X})^{-1}, \end{aligned}$$

sigue que el estimador de $\boldsymbol{\beta}$ en el modelo con datos eliminados, puede ser escrito como:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{(i)} &= (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^\top \mathbf{Y}_{(i)} \\ &= \left\{ \mathbf{I} + \frac{1}{1 - h_{ii}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top \right\} (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{Y} - \mathbf{x}_i Y_i) \\ &= \left\{ \mathbf{I} + \frac{1}{1 - h_{ii}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top \right\} (\hat{\boldsymbol{\beta}} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i Y_i) \\ &= \hat{\boldsymbol{\beta}} - \frac{e_i}{1 - h_{ii}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i. \end{aligned}$$

Este resultado permite definir, por ejemplo,

$$e_{j(i)} = Y_j - \hat{Y}_{j(i)} = Y_j - \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}_{(i)},$$

como el j -ésimo residuo con la i -ésima observación eliminada. Adicionalmente,

$$\begin{aligned} e_{i(i)} &= Y_i - \hat{Y}_{i(i)} = Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{(i)} = Y_i - \mathbf{x}_i^\top \left(\hat{\boldsymbol{\beta}} - \frac{e_i}{1 - h_{ii}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \right) \\ &= e_i + \frac{e_i h_{ii}}{1 - h_{ii}} = \frac{e_i}{1 - h_{ii}}, \end{aligned}$$

es conocido como *residuo eliminado*. Lo que por su vez, permite evaluar el efecto de la i -ésima observación sobre el estimador de σ^2 . Considerare el estimador $s_{(i)}^2$, definido mediante:

$$\text{RSS}_{(i)} = \sum_{j \neq i} (Y_j - \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}_{(i)})^2 = \text{RSS} - \frac{e_i^2}{1 - h_{ii}}.$$

De ahí que, el estimador insesgado de σ^2 cuando removemos la i -ésima observación es dado por:

$$s_{(i)}^2 = \frac{1}{n - p - 1} \left\{ (n - p) s^2 - \frac{e_i^2}{1 - h_{ii}} \right\}$$

RESULTADO 3.21. *Considere el modelo de salto en la media:*

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{d}_i\gamma + \boldsymbol{\epsilon}, \quad (3.8)$$

con $\boldsymbol{\epsilon} \sim \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ y $\mathbf{d}_i = (\mathbf{0}, 1, \mathbf{0})^\top$ un vector de ceros con un 1 en la i -ésima posición. De este modo, el estimador ML de $\boldsymbol{\beta}$ en el modelo (3.7) y (3.8) coinciden.

DEMOSTRACIÓN. El resultado sigue mediante escribir el modelo en (3.8) como

$$\mathbf{Y} = \mathbf{Z}\boldsymbol{\theta} + \boldsymbol{\epsilon}, \quad \mathbf{Z} = (\mathbf{X}, \mathbf{d}_i), \quad \boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \gamma)^\top,$$

y $\hat{\boldsymbol{\theta}} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y}$. Tenemos que

$$\begin{aligned} \mathbf{Z}^\top \mathbf{Z} &= \begin{pmatrix} \mathbf{X}^\top \\ \mathbf{d}_i^\top \end{pmatrix} (\mathbf{X}, \mathbf{d}_i) = \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{d}_i \\ \mathbf{d}_i^\top \mathbf{X} & \mathbf{d}_i^\top \mathbf{d}_i \end{pmatrix} = \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{x}_i \\ \mathbf{x}_i^\top & 1 \end{pmatrix}, \\ \mathbf{Z}^\top \mathbf{Y} &= \begin{pmatrix} \mathbf{X}^\top \\ \mathbf{d}_i^\top \end{pmatrix} \mathbf{Y} = \begin{pmatrix} \mathbf{X}^\top \mathbf{Y} \\ Y_i \end{pmatrix}. \end{aligned}$$

Sabemos que

$$(\mathbf{Z}^\top \mathbf{Z})^{-1} = \begin{pmatrix} (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{1}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} & -\frac{1}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \\ -\frac{1}{1-h_{ii}} \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} & \frac{1}{1-h_{ii}} \end{pmatrix},$$

luego

$$\begin{aligned} \hat{\beta}_* &= \left\{ (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{1}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \right\} \mathbf{X}^\top \mathbf{Y} - \frac{Y_i}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \\ &= \left\{ \mathbf{I} + \frac{1}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^\top \right\} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - \frac{Y_i}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i \\ &= \hat{\beta} + \frac{\hat{Y}_i}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i - \frac{Y_i}{1-h_{ii}}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i = \hat{\beta}_{(i)}, \end{aligned}$$

lo que termina la prueba. \square

Basado en el elipsoide de confianza del $100(1-\alpha)\%$ para β ,

$$\frac{(\beta - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\beta - \hat{\beta})}{ps^2} \leq F_{p,n-p}(1-\alpha).$$

[Cook \(1977\)](#) propuso determinar la influencia de la i -ésima observación, usando

$$D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\beta}_{(i)} - \hat{\beta})}{ps^2} = \frac{r_i^2}{p} \left(\frac{h_i}{1-h_i} \right),$$

para $i = 1, \dots, n$, y recomendó comparar D_i con algún percentil de la distribución $F_{p,n-p}(\alpha)$ con $\alpha = 0.10$. Aunque otra alternativa que puede ser más razonable es usar $\alpha = 0.50$. Adicionalmente se ha sugerido que $D_i > 1$ es un *indicador de observaciones influyentes*.

Se han introducido diversas medidas de diagnóstico, por ejemplo, [Welsch y Kuh \(1977\)](#) propusieron medir el impacto en la i -ésima observación sobre el valor predicho como

$$\text{DFFIT}_i = \hat{Y}_i - \hat{Y}_{i(i)} = \mathbf{x}_i^\top (\hat{\beta} - \hat{\beta}_{(i)}) = \frac{h_{ii} e_i}{1-h_{ii}},$$

o bien, utilizando su versión estandarizada

$$\begin{aligned} \text{DFFITS}_i &= \frac{\hat{Y}_i - \hat{Y}_{i(i)}}{s_{(i)} \sqrt{h_{ii}}} = \left(\frac{h_{ii}}{1-h_{ii}} \right)^{1/2} \frac{e_i}{s_{(i)} \sqrt{1-h_{ii}}} \\ &= \left(\frac{h_{ii}}{1-h_{ii}} \right)^{1/2} t_i. \end{aligned}$$

Sobre este tipo de medidas, [Belsley et al. \(1980\)](#) sugieren poner especial atención en aquellos casos donde $\text{DFFITS}_i > 2\sqrt{p/n}$.

OBSERVACIÓN 3.22. En ocasiones esta medida es conocida como *distancia Welsch-Kuh*.

Por otro lado, [Atkinson \(1981\)](#) sugirió usar una versión modificada de la distancia de Cook, como

$$\begin{aligned} AK_i &= \sqrt{\left(\frac{n-p}{p}\right)\left(\frac{h_{ii}}{1-h_{ii}}\right)} \left| \frac{e_i}{s_{(i)}\sqrt{1-h_{ii}}} \right| \\ &= \sqrt{\left(\frac{n-p}{p}\right)\left(\frac{h_{ii}}{1-h_{ii}}\right)} |t_i| \\ &= \sqrt{\frac{n-p}{p}} |\text{DFFITS}_i|. \end{aligned}$$

Cuando $h_{ii} = p/n, \forall i$, tenemos $AK_i = |t_i|$ debido a esto se recomienda hacer el gráfico de AK_i vs. $|t_i|$. Además podemos identificar la i -ésima observación como influyente si $AK_i > 2$.

OBSERVACIÓN 3.23. AK_i puede considerarse como una medida de influencia conjunta sobre $\hat{\beta}$ y s^2 simultáneamente.

Como un intento de fundamentar las técnicas para llevar a cabo diagnóstico por eliminación de casos, [Cook y Weisberg \(1980\)](#) propusieron considerar medidas generales de influencia, basadas en la función de influencia empírica, lo que los llevó a definir:

$$D_i(\mathbf{M}, c) = \frac{(\hat{\beta} - \hat{\beta}_{(i)})^\top \mathbf{M} (\hat{\beta} - \hat{\beta}_{(i)})}{c}, \quad i = 1, \dots, n,$$

donde \mathbf{M} es matriz definida positiva $p \times p$ y $c > 0$ es un factor de escala. En efecto, podemos escribir diversas medidas bajo esta perspectiva, considere la siguiente tabla:

\mathbf{M}	c	Medida	Referencia
$\mathbf{X}^\top \mathbf{X}$	ps^2	D_i	Cook (1977)
$\mathbf{X}^\top \mathbf{X}$	$ps_{(i)}^2$	$(\text{DFFITS}_i)^2$	Welsch y Kuh (1977)
$\mathbf{X}^\top \mathbf{X}$	$(n-1)^2 ps_{(i)}^2 / (n-p)$	AK_i	Atkinson (1981)

Es posible, centrar nuestra atención en diversos aspectos de la modelación, no solamente en el efecto sobre los estimadores. Por ejemplo, resulta interesante comparar $\text{Cov}(\hat{\beta}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$ con la matriz de covarianza que resulta de eliminar el i -ésimo caso. Esto llevó a [Belsley et al. \(1980\)](#) a definir,

$$\begin{aligned} COVRATIO_i &= \frac{\det\{s_{(i)}^2(\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1}\}}{\det\{s^2(\mathbf{X}^\top \mathbf{X})^{-1}\}} = \left(\frac{s_{(i)}^2}{s^2}\right)^p \frac{\det(\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1}}{\det(\mathbf{X}^\top \mathbf{X})^{-1}} \\ &= \frac{1}{1-h_i} \left(\frac{n-p-r_i^2}{n-p-1}\right)^p. \end{aligned}$$

Además, se ha planteado el siguiente punto de corte $|COVRATIO_i - 1| > 3p/n$.

Debemos destacar que [Belsley et al. \(1980\)](#) y [Velleman y Welsch \(1981\)](#) han discutido *estrategias para amenizar el cálculo* de estas medidas de influencia. Para modelos de regresión lineal algunas de estas medidas han sido implementadas en software estadístico tales como SAS, SPSS o S-Plus/R.

En particular, R (o S-Plus) disponen de las funciones `lm.influence` y `ls.diag` asociadas con las funciones `lm` (o `glm`) y `lsfit`, respectivamente. Estas cantidades pueden ser obtenidas de *forma eficiente* usando la descomposición QR o SVD.

3.6. Procedimientos para estimación robusta

En esta sección presentamos dos métodos para obtener estimadores que permiten atenuar el efecto de outliers o observaciones atípicas. El primero es conocido como M -estimadores, mientras que el segundo se basa en considerar distribuciones con colas más pesadas que la normal.

3.6.1. M -estimadores. Para introducir ideas considere X_1, \dots, X_n variables aleatorias IID. Sabemos que la media muestral \bar{X} es solución del problema,

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

o análogamente,

$$\sum_{i=1}^n (x_i - \theta) = 0.$$

Mientras que la mediana, denotada como $\text{me}(\mathbf{x})$, que es *robusta contra outliers*, es solución del problema

$$\min_{\theta} \sum_{i=1}^n |x_i - \theta|,$$

es decir $\text{me}(\mathbf{x})$ debe satisfacer la ecuación

$$\sum_{i=1}^n \{(-1)I_{(-\infty, 0)}(x_i - \theta) + I_{(0, \infty)}(x_i - \theta)\} = 0.$$

Sabemos que para una muestra aleatoria X_1, \dots, X_n desde el modelo estadístico $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$ con función de densidad $f(x; \theta)$, el MLE $\hat{\theta}_{\text{ML}}$ corresponde al minimizador del negativo de la log-verosimilitud

$$\min_{\theta} \left\{ - \sum_{i=1}^n \log f(x_i; \theta) \right\},$$

y en el caso de que $\ell(\theta) = \sum_{i=1}^n \log f(x_i; \theta)$ sea diferenciable, $\hat{\theta}_{\text{ML}}$ debe ser solución de:

$$\sum_{i=1}^n \frac{\partial \log f(x_i; \theta)}{\partial \theta} = 0,$$

lo que lleva a la siguiente definición.

DEFINICIÓN 3.24 (M -estimador). Para X_1, \dots, X_n variables aleatorias IID. El M -estimador $\hat{\theta}_{\text{M}}$ con respecto a la función $\psi : \mathbb{R} \times \Theta \rightarrow \mathbb{R}$ se define como la solución de la ecuación

$$\sum_{i=1}^n \psi(X_i; \hat{\theta}_{\text{M}}) = 0. \quad (3.9)$$

Usualmente el estimador definido por la Ecuación (3.9) corresponde a la solución del problema de optimización

$$\min_{\theta} \sum_{i=1}^n \rho(x_i; \theta),$$

si ρ es diferenciable, entonces

$$\psi(x; \theta) = c \frac{\partial \rho(x; \theta)}{\partial \theta},$$

para alguna constante c .

OBSERVACIÓN 3.25. Para simplificar la notación podemos hacer

$$\rho(x; \theta) = \tilde{\rho}(x - \theta), \quad \psi(x; \theta) = \tilde{\psi}(x - \theta).$$

EJEMPLO 3.26 (Estimador LS). Sea $z = x - \theta$. Las funciones $\tilde{\rho}(z) = z^2$ y $\tilde{\psi}(z) = z$ llevan a la media muestral.

EJEMPLO 3.27 (Estimador LAD³). La mediana es un M -estimador con $\tilde{\rho}(z) = |z|$, y

$$\tilde{\psi}(z) = \begin{cases} -1, & z < 0, \\ 0, & z = 0, \\ 1, & z > 0. \end{cases}$$

EJEMPLO 3.28 (Media recortada). La primera propuesta de Huber (1981) para reducir la influencia de outliers es:

$$\tilde{\rho}(z) = \begin{cases} z^2, & |z| \leq k, \\ k^2, & |z| > k, \end{cases}$$

donde k es una *constante de tuning*, y

$$\tilde{\psi}(z) = \begin{cases} z, & |z| \leq k, \\ 0, & |z| > k. \end{cases}$$

EJEMPLO 3.29 (Media Winsorizada). Huber (1981) propuso un compromiso entre la media y la mediana, como:

$$\tilde{\rho}(z) = \begin{cases} \frac{1}{2}z^2, & |z| \leq k, \\ k|z| - \frac{1}{2}k^2, & |z| > k. \end{cases}$$

El estimador $\hat{\theta}_M$ es solución de:

$$\sum_{i=1}^n \tilde{\psi}(x_i - \theta) = 0,$$

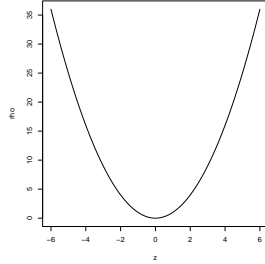
donde

$$\tilde{\psi}(z) = \begin{cases} -k, & z < -k, \\ z, & |z| \leq k, \\ k, & z > k. \end{cases}$$

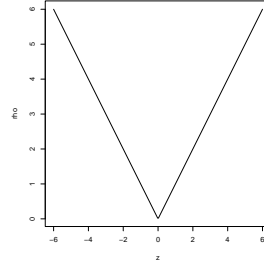
Debemos destacar que, para $k \rightarrow 0$ obtenemos la mediana, cuando $k \rightarrow \infty$ lleva a la media, mientras que $k = 1.345$ tiene 95% de eficiencia bajo normalidad.

³Estimadores *mínimo desvío absoluto* (LAD) corresponden a estimadores L_1 .

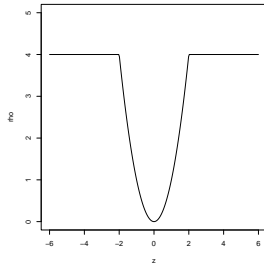
En los gráficos a continuación se presenta las funciones $\rho(\cdot)$ para cada uno de los ejemplos anteriores,



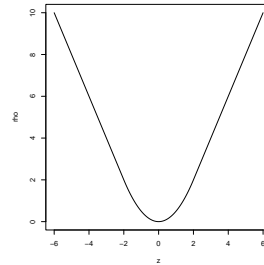
(a) media



(b) mediana



(c) media recortada



(d) media winsorizada

Sabemos que el estimador LS en regresión es solución del problema de estimación

$$\min_{\beta} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \beta)^2, \quad (3.10)$$

y, bajo normalidad, el MLE de β minimiza la función

$$-\sum_{i=1}^n \log f(Y_i; \beta) = \frac{n}{2} \log 2\pi\sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \beta)^2. \quad (3.11)$$

Para obtener estimadores robustos [Huber \(1981\)](#) sugirió substituir el negativo de la función de log-verosimilitud en (3.11) por una función que permita disminuir el efecto de outliers. De este modo [Huber \(1981\)](#) propuso obtener estimadores tipo-ML, conocidos como M -estimadores resolviendo el problema

$$\min_{\beta} \sum_{i=1}^n \rho(Y_i - \mathbf{x}_i^\top \beta). \quad (3.12)$$

Es usual incorporar el parámetro de escala σ^2 definiendo,

$$z_i = \frac{Y_i - \mathbf{x}_i^\top \beta}{\sigma}, \quad i = 1, \dots, n,$$

lo que lleva a considerar la función objetivo:

$$Q_\rho(\beta) = \sum_{i=1}^n \rho(z_i).$$

Es fácil notar que,

$$\frac{\partial Q_\rho(\boldsymbol{\beta})}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial \rho(z_i)}{\partial z_i} \frac{\partial z_i}{\partial \beta_j} = \sum_{i=1}^n W(z_i) z_i \frac{\partial z_i}{\partial \beta_j}, \quad W(z_i) = \frac{1}{z_i} \frac{\partial \rho(z_i)}{\partial z_i}.$$

Tenemos que $\psi(z) = \partial \rho(z)/\partial z$, así $W(z) = \psi(z)/z$. Por tanto el M -estimador de $\boldsymbol{\beta}$ es solución del sistema de ecuaciones

$$\sum_{i=1}^n \psi\left(\frac{Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}}{\sigma}\right) \frac{x_{ij}}{\sigma} = 0, \quad j = 1, \dots, p,$$

que puede ser escrito en forma compacta como

$$\frac{1}{\sigma^2} \sum_{i=1}^n W_i (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}) \mathbf{x}_i = \mathbf{0} \quad (3.13)$$

Sea $\mathbf{W} = \text{diag}(W_1, \dots, W_n)$ de este modo podemos escribir (3.13) como:

$$\mathbf{X}^\top \mathbf{W} (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) = \mathbf{0}.$$

Sin embargo, debemos resaltar que los pesos W_i dependen de z_i , los que a su vez dependen de $\boldsymbol{\beta}$ y por tanto se requiere métodos iterativos para obtener el M -estimador, $\hat{\boldsymbol{\beta}}_M$.

Usando una estimación inicial $\boldsymbol{\beta}^{(0)}$ podemos considerar el siguiente esquema iterativo

$$\boldsymbol{\beta}^{(r+1)} = (\mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{Y}, \quad (3.14)$$

con

$$\mathbf{W}^{(r)} = \text{diag}(W_1^{(r)}, \dots, W_n^{(r)}), \quad W_i^{(r)} = \psi(e_i^{(r)}/\sigma)/(e_i^{(r)}/\sigma),$$

y

$$\mathbf{e}^{(r)} = \mathbf{Y} - \mathbf{X} \boldsymbol{\beta}^{(r)}.$$

Desde el punto de vista computacional es preferible usar

$$\boldsymbol{\beta}^{(r+1)} = \boldsymbol{\beta}^{(r)} + \mathbf{p}_r,$$

con

$$\mathbf{p}_r = (\mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{e}^{(r)}.$$

OBSERVACIÓN 3.30. El procedimiento delineado en (3.14) es conocido como *mínimos cuadrados iterativamente ponderados (IRLS)*. Debemos destacar que la convergencia del proceso iterativo en (3.14) sólo es garantizada para funciones ρ *convexas* así como para funciones de *redescenso*. Además, una serie de autores han propuesto métodos refinados para obtener M -estimadores basados en IRLS (ver, por ejemplo [O'Leary, 1990](#)).

Existe una gran variedad de funciones $\rho(\cdot)$ o $\psi(\cdot)$ para definir M -estimadores, por ejemplo:

- *Tukey's biweight*

$$\psi(z) = z[1 - (t/k)]_+^2, \quad k = 4.685.$$

- *Hampel's ψ*

$$\psi(z) = \text{sign}(z) \begin{cases} |z|, & 0 < |z| \leq a, \\ a, & a < |z| \leq b, \\ a(c - |z|)/(c - b), & b < |z| \leq c, \\ 0, & c < |z|, \end{cases}$$

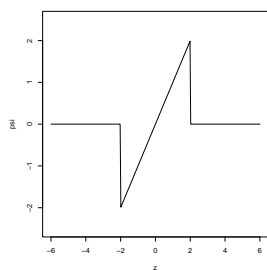
con $a = 1.645$, $b = 3$ y $c = 6.5$.

- *Función seno de Andrews*

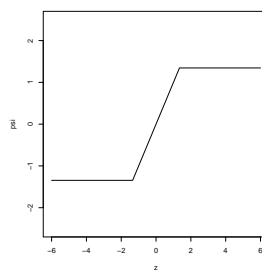
$$\psi(z) = \begin{cases} \sin(z/a), & |z| \leq \pi a, \\ 0, & |z| > \pi a, \end{cases}$$

donde $a = 1.339$.

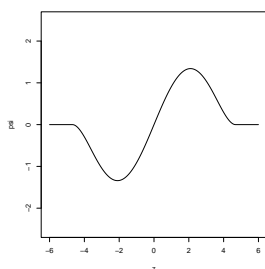
Resulta instructivo visualizar el comportamiento de algunas funciones $\psi(\cdot)$ seleccionadas en la figura a continuación,



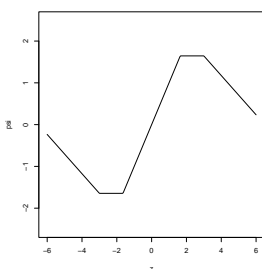
(a) media recortada



(b) Huber



(c) Tukey bisquare



(d) Hampel

Aunque es de interés primario la estimación de los coeficientes de regresión, se ha sugerido usar el siguiente estimador robusto para σ ,

$$\hat{\sigma}_{\text{rob}} = \frac{\text{MAD}(\mathbf{e})}{0.6745},$$

con

$$\text{MAD}(\mathbf{e}) = \text{me}(|\mathbf{e} - \text{me}(\mathbf{e})|),$$

que es conocido como *desviación mediana absoluta* de los residuos minimos cuadrados $\mathbf{e} = \mathbf{Y} - \mathbf{X}\hat{\beta}_{\text{LS}}$.

Análogamente a β es posible obtener un M -estimador para σ resolviendo

$$\frac{1}{n-p} \sum_{i=1}^n \chi\left(\frac{Y_i - \mathbf{x}_i^\top \beta}{\sigma}\right) = \delta,$$

para δ una constante positiva.

La fundamentación de las propiedades asintóticas de la clase de M -estimadores de β escapan al ámbito en el que se desarrolla la asignatura. Sin embargo, podemos proveer un estimador de la matriz de covarianza de $\hat{\beta}_M$ considerando,

$$\kappa = 1 + \frac{p}{n} \frac{\text{var}(\psi')}{\{\mathbb{E}(\psi')\}^2},$$

que es evaluado en la distribución de los errores ϵ (en la práctica deben ser estimados desde los residuos). Entonces la matriz de covarianza asintótica de $\hat{\beta}_M$ es dada por (ver [Huber, 1981](#))

$$\kappa^2 \frac{\sum_i \psi^2(e_i)/(n-p)}{[\sum_i \psi'(e_i)/n]^2} (\mathbf{X}^\top \mathbf{X})^{-1}.$$

Una subclase interesante de M -estimadores, corresponde a los problemas de estimación norma- L_p , los que son definidos como la solución del problema:

$$\min_{\beta} \left(\sum_{i=1}^n |Y_i - \mathbf{x}_i^\top \beta|^p \right)^{1/p}, \quad p \geq 1. \quad (3.15)$$

En efecto, basta considerar

$$\rho(z) = |z|^p, \quad \psi(z) = |z|^{p-2}, \quad W(z) = |z|^{p-2}.$$

OBSERVACIÓN 3.31. Considere Y_1, \dots, Y_n variables aleatorias provenientes de la función de densidad

$$f_p(y) = c \exp(-|y|^p), \quad p \geq 1,$$

con c una constante de normalización. Esto permite caracterizar la estimación de norma- L_p como estimación ML basado en la densidad $f_p(y)$. Además, se debe destacar que la *distribución Laplace* es obtenida para $p = 1$, mientras que la *distribución normal* es recuperada para $p = 2$.

Para $1 < p < 2$ podemos usar IRLS como un método para aproximar la solución del problema en (3.15). En efecto, [Osborne \(1985\)](#) reestableció el problema de estimación norma- L_p como:

$$\min_{\beta} Q_p(\beta), \quad Q_p(\beta) = \sum_{i=1}^n |\epsilon_i|^p = \sum_{i=1}^n |\epsilon_i|^{p-2} \epsilon_i^2,$$

que puede ser interpretado como un problema de mínimos cuadrados ponderados,

$$\min_{\beta} \|\mathbf{W}^{(p-2)/2} (\mathbf{Y} - \mathbf{X}\beta)\|_2^2, \quad \mathbf{W} = \text{diag}(|\epsilon_1|, \dots, |\epsilon_n|).$$

Esto lleva al siguiente algoritmo.

Algoritmo 5: IRLS para estimación L_p .**Entrada:** Datos \mathbf{X} , \mathbf{y} , $p \in [1, 2)$ y estimación inicial $\boldsymbol{\beta}^{(0)}$.**Salida :** Aproximación del estimador L_p , $\hat{\boldsymbol{\beta}}$.

```

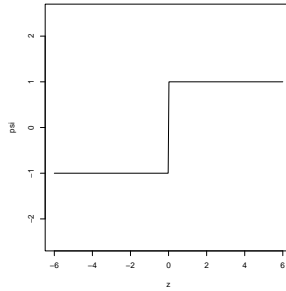
1 begin
2   for  $r = 0, 1, 2, \dots$  do
3      $\mathbf{e}^{(r)} = \mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^{(r)}$ 
4      $\mathbf{W}^{(r)} = \text{diag}(|e_1^{(r)}|^{(p-2)/2}, \dots, |e_n^{(r)}|^{(p-2)/2})$ 
5     Resolver  $\mathbf{p}_r$  desde
6        $\min_{\mathbf{p}_r} \|\mathbf{W}^{(r)}(\mathbf{e}^{(r)} - \mathbf{X}\mathbf{p}_r)\|_2^2$ 
7     Hacer
8        $\boldsymbol{\beta}^{(r+1)} = \boldsymbol{\beta}^{(r)} + \mathbf{p}_r$ 
9   end
10  return  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^{(*)}$ 
11 end

```

Debemos destacar que, para $p = 1$ tenemos

$$\psi(z) = z/|z|.$$

con



De este modo, el procedimiento de estimación norma- L_1 es robusto, con pesos

$$W_i = 1/|e_i|, \quad i = 1, \dots, n.$$

Schlossmacher (1973) fue uno de los primeros autores en proponer el uso de IRLS para estimación norma- L_1 en regresión basado en la ecuación

$$\sum_{i=1}^n \frac{Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}}{|Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}|} \mathbf{x}_i = \mathbf{0}.$$

Sin embargo este procedimiento ha sido fuertemente criticado por una gran cantidad de autores (ver discusión en Capítulo 4 de Björck, 1996), por no ser capaz de identificar las *observaciones básicas* que definen el estimador.

Una perspectiva diferente fue adoptada por Charnes et al. (1955) quienes mostraron que problema de regresión L_1 ,

$$\min_{\boldsymbol{\beta}} Q_1(\boldsymbol{\beta}), \quad Q_1(\boldsymbol{\beta}) = \sum_{i=1}^n |Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}|,$$

es equivalente a resolver el siguiente problema de *programación lineal* (LP)

$$\begin{aligned} & \min \sum_{i=1}^n (\epsilon_i^+ + \epsilon_i^-), \\ \text{sujeto a: } & \epsilon_i^+ \geq 0, \quad \epsilon_i^- \geq 0, \\ & \epsilon_i^+ - \epsilon_i^- = \epsilon_i = Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}, \end{aligned}$$

para $i = 1, \dots, n$, con ϵ_i^+ y ϵ_i^- variables no negativas.

OBSERVACIÓN 3.32. Barrodale y Roberts (1973, 1974) presentaron un algoritmo de propósito especial para resolver el problema de regresión L_1 , modificando el método simplex y la estructura de datos requerida.

3.6.2. Estimación usando distribuciones con colas pesadas. Frecuentemente se ha sugerido substituir la distribución normal por distribuciones con colas más pesadas como un mecanismo para la *acomodación de outliers*. En efecto, en la Observación 3.31 se indicó que usar la distribución Laplace, permite obtener estimadores L_1 en regresión, lo que corresponde a un método robusto contra outliers. En esta sección sin embargo, confiaremos en la simpleza y elegancia del procedimiento de estimación por máxima verosimilitud. En esta sección nos enfocaremos en abordar la estimación de los parámetros de regresión usando la clase de distribuciones de contornos elípticos, discutida en la Sección 1.6.

Debemos resaltar que *diferentemente* al caso de la distribución normal, en el caso general de la familia elíptica, podemos tener los siguientes enfoques:

- (a) *Modelo dependiente:* Supondremos Y_1, \dots, Y_n tal que su *densidad conjunta* $\mathbf{Y} = (Y_1, \dots, Y_n)^\top \sim \text{EC}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n; g)$, sigue una distribución de contornos elípticos.⁴
- (b) *Modelo independiente:* Considere Y_1, \dots, Y_n variables aleatorias *independientes* cada una con distribución $\text{EC}_1(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; g)$.

Debido a esta consideración a continuación nos enfocaremos en el *modelo independiente*,

$$Y_i \stackrel{\text{ind}}{\sim} \text{EC}_1(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; g), \quad i = 1, \dots, n,$$

donde

$$f(y_i; \boldsymbol{\theta}) = \frac{1}{\sigma} g((Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 / \sigma^2),$$

con $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$, y la función generadora de densidades g debe satisfacer:

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

OBSERVACIÓN 3.33. Como se mencionó en la Sección 1.6, la función generadora de densidades $g(\cdot)$ puede depender de *parámetros de forma*, los que controlan el grado de curtosis de la distribución (es decir, que tan *pesada* es la cola de la distribución). A continuación asumiremos que tales parámetros son conocidos, o equivalentemente que la función $g(\cdot)$ es conocida.

⁴El estimador de $\boldsymbol{\beta}$ en el modelo dependiente es equivalente al LSE y por tanto NO es robusto.

Obtenemos los estimadores máximo verosímiles de β y σ^2 , asumiendo $g(\cdot)$ conocido, mediante maximizar la función de log-verosimilitud:

$$\ell(\theta) = -\frac{n}{2} \log \sigma^2 + \sum_{i=1}^n \log g(u_i),$$

donde $u_i = (Y_i - \mathbf{x}_i^\top \beta)^2 / \sigma^2$, para $i = 1, \dots, n$. Diferenciando $\ell(\theta)$ con relación a β , obtenemos

$$\begin{aligned} d_\beta \ell(\theta) &= \sum_{i=1}^n d_\beta \log g(u_i) = \sum_{i=1}^n \frac{g'(u_i)}{g(u_i)} d_\beta u_i \\ &= -\frac{2}{\sigma^2} \sum_{i=1}^n \frac{g'(u_i)}{g(u_i)} (Y_i - \mathbf{x}_i^\top \beta) \mathbf{x}_i^\top d\beta \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n W_i(\theta) (Y_i - \mathbf{x}_i^\top \beta) \mathbf{x}_i^\top d\beta \end{aligned}$$

donde $W_i(\theta) = -2g'(u_i)/g(u_i)$, para $i = 1, \dots, n$. Análogamente, diferenciando $\ell(\theta)$ con relación a σ^2 , tenemos

$$\begin{aligned} d_{\sigma^2} \ell(\theta) &= -\frac{n}{2\sigma^2} d\sigma^2 + \sum_{i=1}^n \frac{g'(u_i)}{g(u_i)} d_{\sigma^2} u_i \\ &= -\frac{n}{2\sigma^2} d\sigma^2 - \frac{1}{\sigma^4} \sum_{i=1}^n \frac{g'(u_i)}{g(u_i)} (Y_i - \mathbf{x}_i^\top \beta)^2 d\sigma^2 \\ &= -\frac{n}{2\sigma^2} d\sigma^2 + \frac{1}{2\sigma^4} \sum_{i=1}^n W_i(\theta) (Y_i - \mathbf{x}_i^\top \beta)^2 d\sigma^2 \\ &= -\frac{n}{2\sigma^2} d\sigma^2 + \frac{1}{2\sigma^4} (\mathbf{Y} - \mathbf{X}\beta)^\top \mathbf{W} (\mathbf{Y} - \mathbf{X}\beta) d\sigma^2 \end{aligned}$$

donde

$$\mathbf{W} = \text{diag} (W_1(\theta), \dots, W_n(\theta)).$$

La condición de primer orden, lleva al siguiente sistema de ecuaciones:

$$\begin{aligned} \mathbf{X}^\top \widehat{\mathbf{W}} (\mathbf{Y} - \mathbf{X}\widehat{\beta}) &= \mathbf{0}, \\ \widehat{\sigma}^2 &= \frac{1}{n} (\mathbf{Y} - \mathbf{X}\widehat{\beta})^\top \widehat{\mathbf{W}} (\mathbf{Y} - \mathbf{X}\widehat{\beta}), \end{aligned}$$

que *no tiene solución en forma explícita* y por tanto se requiere el uso de métodos iterativos.

Por ejemplo, usando una estimación inicial $\theta = \theta^{(k)}$, actualizamos las estimaciones para β y σ^2 , como:

$$\begin{aligned} \beta^{(k+1)} &= (\mathbf{X}^\top \mathbf{W}^{(k)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(k)} \mathbf{Y}, \\ \sigma^{2(k+1)} &= \frac{1}{n} (\mathbf{Y} - \mathbf{X}\beta^{(k)})^\top \mathbf{W}^{(k)} (\mathbf{Y} - \mathbf{X}\beta^{(k)}), \end{aligned}$$

a la convergencia del algoritmo, hacemos $(\widehat{\beta}, \widehat{\sigma}^2)$.

EJEMPLO 3.34. Las funciones de pesos para algunas distribuciones elípticas seleccionadas son dadas a continuación:

- *Normal*: $Y_1, \dots, Y_n \stackrel{\text{ind}}{\sim} \mathbf{N}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2)$, tenemos:

$$W_i(\boldsymbol{\theta}) = 1, \quad i = 1, \dots, n.$$

- *t-Student*: $Y_1, \dots, Y_n \stackrel{\text{ind}}{\sim} t(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; \nu)$, $\nu > 0$,

$$W_i(\boldsymbol{\theta}) = \frac{\nu + 1}{\nu + u_i}, \quad i = 1, \dots, n,$$

la distribución $\text{Cauchy}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2)$, es obtenida para $\nu = 1$.

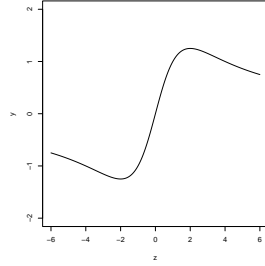
- *Normal contaminada*: $Y_1, \dots, Y_n \stackrel{\text{ind}}{\sim} \text{CN}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; \epsilon, \gamma)$, con $\epsilon \in [0, 1)$, $\gamma > 0$,

$$W_i(\boldsymbol{\theta}) = \frac{(1 - \epsilon) \exp(-u/2) + \epsilon \gamma^{-3/2} \exp(-u/(2\gamma))}{(1 - \epsilon) \exp(-u/2) + \epsilon \gamma^{-1/2} \exp(-u/(2\gamma))}, \quad i = 1, \dots, n.$$

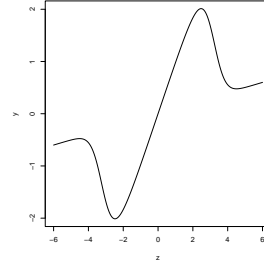
- *Exponencial Potencia*: $Y_1, \dots, Y_n \stackrel{\text{ind}}{\sim} \text{PE}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; \lambda)$, $\lambda > 0$,

$$W_i(\boldsymbol{\theta}) = \lambda u_i^{\lambda-1}, \quad i = 1, \dots, n.$$

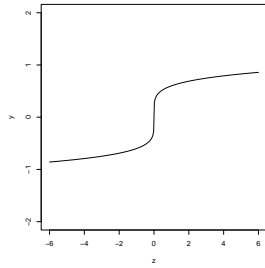
Es interesante notar el comportamiento de las funciones de influencia para cada una de estas distribuciones, considere el siguiente gráfico:



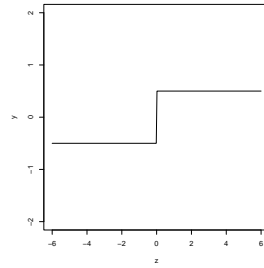
(a) *t* de Student, $\nu = 4$



(b) CN, $\epsilon = 0.05, \gamma = 10$



(c) PE, $\lambda = 0.6$



(d) PE, $\lambda = 0.5$

La clase de distribuciones elípticas contiene distribuciones con colas más pesadas y también más livianas que la normal. A continuación nos enfocaremos en una subclase conocida como *mezclas de escala normal* (Andrews y Mallows, 1974), clase que tiene una interesante interpretación que será explotada para la estimación de parámetros.

Sea Y_1, \dots, Y_n variable aleatorias independientes con distribución $\text{SMN}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; \mathbf{H})$ cada una con densidad

$$f(y_i; \boldsymbol{\theta}) = (2\pi\sigma)^{-1/2} \int_0^\infty \omega^{1/2} \exp(-\omega u_i/2) d\mathbf{H}(\omega),$$

donde $u_i = (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 / \sigma^2$, $i = 1, \dots, n$.

OBSERVACIÓN 3.35. Sabemos que, una variable aleatoria $Y_i \sim \text{SMN}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; \mathbf{H})$, admite la representación:

$$Y_i | W = \omega \sim \mathbf{N}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2 / \omega), \quad W \sim \mathbf{H}(\boldsymbol{\delta}),$$

lo que permite abordar la estimación ML usando el *algoritmo EM* (Dempster, Laird y Rubin, 1977).

El algoritmo EM permite el cálculo iterativo de *estimadores ML* en modelos con *datos incompletos*. De este modo, requiere de una *formulación de datos aumentados*. Su principal ventaja es que reemplaza una optimización “*compleja*” (asociada a la estimación ML) por una serie de maximizaciones “*simples*”. A continuación damos una breve descripción de este procedimiento, para una revisión de las propiedades de este algoritmo y una serie de consideraciones prácticas consulte McLachlan y Krishnan (2008).

3.6.2.1. *Algoritmo EM (Esperanza-Maximización)*. Sea \mathbf{Y}_{obs} vector de *datos observados* con función de densidad $f(\mathbf{y}_{\text{obs}}; \boldsymbol{\theta})$. El objetivo es aumentar los datos observados \mathbf{Y}_{obs} con variables latentes \mathbf{Y}_{mis} también conocidos como *datos perdidos*. Esto es, consideramos el vector de *datos completos*

$$\mathbf{Y}_{\text{com}} = (\mathbf{Y}_{\text{obs}}^\top, \mathbf{Y}_{\text{mis}}^\top)^\top,$$

tal que la densidad $f(\mathbf{y}_{\text{com}}; \boldsymbol{\theta})$ sea simple. El algoritmo EM es útil cuando la función de log-verosimilitud

$$\begin{aligned} \ell_o(\boldsymbol{\theta}; \mathbf{Y}_{\text{obs}}) &= \log f(\mathbf{y}_{\text{obs}}; \boldsymbol{\theta}) \\ &= \log \int f(\mathbf{y}_{\text{com}}; \boldsymbol{\theta}) d\mathbf{y}_{\text{mis}}, \end{aligned}$$

es difícil de maximizar directamente. El objetivo del algoritmo EM es realizar la estimación ML iterativamente basandose en la *log-verosimilitud de datos completos*:

$$\ell_c(\boldsymbol{\theta}; \mathbf{Y}_{\text{com}}) = \log f(\mathbf{y}_{\text{com}}; \boldsymbol{\theta}).$$

El algoritmo EM permite obtener los MLE en *problemas con datos incompletos* por medio de las etapas:

Paso E: para $\boldsymbol{\theta}^{(k)}$ estimación de $\boldsymbol{\theta}$ en la k -ésima iteración, calcular la Q -función,

$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)}) &= \mathbb{E}\{\ell_c(\boldsymbol{\theta}; \mathbf{Y}_{\text{com}}) | \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}^{(k)}\} \\ &= \int \ell_c(\boldsymbol{\theta}; \mathbf{Y}_{\text{com}}) f(\mathbf{y}_{\text{mis}} | \mathbf{y}_{\text{obs}}, \boldsymbol{\theta}^{(k)}) d\mathbf{y}_{\text{mis}}. \end{aligned}$$

Paso M: determinar $\boldsymbol{\theta}^{(k+1)}$ como

$$\boldsymbol{\theta}^{(k+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)}).$$

Adicionalmente, Dempster, Laird y Rubin (1977) definieron un *algoritmo EM generalizado (GEM)*, mediante la siguiente modificación del Paso M:

Paso M*: seleccionar $\boldsymbol{\theta}^{(k+1)}$ satisfaciendo,

$$Q(\boldsymbol{\theta}^{(k+1)}; \hat{\boldsymbol{\theta}}^{(k)}) > Q(\boldsymbol{\theta}^{(k)}; \boldsymbol{\theta}^{(k)}).$$

OBSERVACIÓN 3.36. Por ejemplo, podemos considerar *un único* paso Newton en la optimización de $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)})$ para definir un GEM.

Algunas propiedades del algoritmo EM se presentan en los siguientes resultados

TEOREMA 3.37. *Todo algoritmo EM o GEM incrementa la log-verosimilitud de datos observados $\ell_o(\boldsymbol{\theta}; \mathbf{Y}_{\text{obs}})$ en cada iteración, esto es,*

$$\ell_o(\boldsymbol{\theta}^{(k+1)}; \mathbf{Y}_{\text{obs}}) \geq \ell_o(\boldsymbol{\theta}^{(k)}; \mathbf{Y}_{\text{obs}}).$$

DEMOSTRACIÓN. Ver [Dempster, Laird y Rubin \(1977\)](#). \square

TEOREMA 3.38 (Convergencia del Algoritmo EM). *Bajo condiciones suaves, la secuencia $\{\boldsymbol{\theta}^{(k)}\}_{k \geq 0}$ generada por el algoritmo EM (GEM). Converge a un punto estacionario de $\ell_o(\boldsymbol{\theta}; \mathbf{Y}_{\text{obs}})$.*

DEMOSTRACIÓN. Ver [Wu \(1983\)](#). \square

Debemos destacar una serie de características relevantes del algoritmo EM, a saber:

- Frecuentemente el algoritmo EM es *simple*, de bajo costo computacional y numéricamente *estable*.
- [Dempster, Laird y Rubin \(1977\)](#) mostraron que el algoritmo EM converge con velocidad lineal, que depende de la *proporción* de información perdida.⁵
- Para modelos con datos aumentados con densidad en la familia exponencial, el algoritmo EM se reduce a *actualizar* las estadísticas suficientes.
- Errores estándar pueden ser obtenidos por cálculo directo, diferenciación numérica o usando el *Principio de Información Perdida* ([Louis, 1982](#)).

Sea Y_1, \dots, Y_n variables aleatorias independientes $\text{SMN}(\mathbf{x}_i \boldsymbol{\beta}, \sigma^2; \mathbf{H})$. Se llevará a cabo la estimación ML usando el algoritmo EM.

De este modo, tenemos el siguiente *modelo jerárquico*:

$$Y_i | W_i = \omega_i \sim \text{N}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2 / \omega_i), \quad W_i \sim \text{H}(\boldsymbol{\delta}), \quad i = 1, \dots, n.$$

En este caso el vector de datos completos es $\mathbf{Y}_{\text{com}} = (\mathbf{Y}^\top, \boldsymbol{\omega}^\top)^\top$, donde

$$\mathbf{Y} = (Y_1, \dots, Y_n)^\top, \quad \boldsymbol{\omega} = (\omega_1, \dots, \omega_n)^\top.$$

En este contexto, \mathbf{Y} corresponde a los datos observados, mientras que $\boldsymbol{\omega}$ serán asumidos como datos perdidos.

Considere una estimación para $\boldsymbol{\theta} = \boldsymbol{\theta}^{(k)}$, entonces

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(k)}) = \text{E}\{\ell_c(\boldsymbol{\theta}; \mathbf{y}_{\text{com}}) | \mathbf{y}; \boldsymbol{\theta}^{(k)}\} = Q_1(\boldsymbol{\beta}, \sigma^2; \boldsymbol{\theta}^{(k)}) + Q_2(\boldsymbol{\delta}; \boldsymbol{\theta}^{(k)}),$$

donde

$$Q_1(\boldsymbol{\beta}, \sigma^2; \boldsymbol{\theta}^{(k)}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n \omega_i^{(k)} (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2,$$

$$Q_2(\boldsymbol{\delta}; \boldsymbol{\theta}^{(k)}) = \text{E}\{\log h^{(n)}(\boldsymbol{\omega}; \boldsymbol{\delta}) | \mathbf{y}; \boldsymbol{\theta}^{(k)}\},$$

⁵puede ser **extremadamente** lento.

con $\omega_i^{(k)} = E(\omega_i | \mathbf{x}_i; \boldsymbol{\theta}^{(k)})$ para $i = 1, \dots, n$. En general, la forma para la esperanza condicional requerida en el Paso-E del algoritmo EM es dada por:

$$E(\omega_i | \mathbf{x}_i; \boldsymbol{\theta}) = \frac{\int_0^\infty \omega_i^{3/2} \exp(-\omega_i u_i/2) dH(\delta)}{\int_0^\infty \omega_i^{1/2} \exp(-\omega_i u_i/2) dH(\delta)},$$

con $u_i = (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 / \sigma^2$, $i = 1, \dots, n$.

A continuación algunos ejemplos de funciones de pesos para miembros en la familia de mezclas de escala normal,

- *t-Student*: $Y_i \stackrel{\text{ind}}{\sim} t(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2, \nu)$, $i = 1, \dots, n$, en cuyo caso

$$E(\omega_i | Y_i; \boldsymbol{\theta}) = \frac{\nu + 1}{\nu + u_i}.$$

- *Slash*: $Y_i \stackrel{\text{ind}}{\sim} \text{Slash}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2, \nu)$, para $i = 1, \dots, n$. De este modo,

$$E(\omega_i | Y_i; \boldsymbol{\theta}) = \left(\frac{2\nu + 1}{u_i} \right) \frac{P_1(\nu + 3/2, u_i/2)}{P_1(\nu + 1/2, u_i/2)},$$

donde

$$P_z(a, b) = \frac{b^a}{\Gamma(a)} \int_0^z t^{a-1} e^{-bt} dt,$$

es la función gama incompleta (regularizada).

- *Exponencial Potencia*: $Y_i \stackrel{\text{ind}}{\sim} \text{PE}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2, \lambda)$, $i = 1, \dots, n$, donde

$$E(\omega_i | Y_i; \boldsymbol{\theta}) = \lambda u_i^{\lambda-1}, \quad u_i \neq 0, \lambda \in (0, 1].$$

Finalmente, el algoritmo EM para obtener los estimadores ML en el modelo

$$Y_i \stackrel{\text{ind}}{\sim} \text{SMN}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2; H), \quad i = 1, \dots, n,$$

adopta la forma:

Paso E: para $\boldsymbol{\theta}^{(k)}$, calcular:

$$\omega_i^{(k)} = E(\omega_i | Y_i; \boldsymbol{\theta}^{(k)}), \quad i = 1, \dots, n.$$

Paso M: actualizar $\boldsymbol{\beta}^{(k+1)}$ y $\sigma^{2(k+1)}$ como:

$$\begin{aligned} \boldsymbol{\beta}^{(k+1)} &= (\mathbf{X}^\top \mathbf{W}^{(k)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(k)} \mathbf{Y}, \\ \sigma^{2(k+1)} &= \frac{1}{n} (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}^{(k)})^\top \mathbf{W}^{(k)} (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}^{(k)}), \end{aligned}$$

donde $\mathbf{W}^{(k)} = \text{diag}(\omega_1(\boldsymbol{\theta}^{(k)}), \dots, \omega_n(\boldsymbol{\theta}^{(k)}))$.

Iteramos entre las etapas E y M hasta alcanzar convergencia, en cuyo caso hacemos $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$ y $\sigma^2 = \hat{\sigma}^2$.

Usando resultados asintóticos tradicionales asociados a la estimación máximo verosímil, podemos notar que

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \mathcal{F}^{-1}(\boldsymbol{\beta})),$$

donde la matriz de información de Fisher asociada a los coeficientes de regresión $\boldsymbol{\beta}$ es dada por (ver por ejemplo, [Lange, Little y Taylor, 1989](#))

$$\mathcal{F}(\boldsymbol{\beta}) = \frac{4\alpha}{\phi} \mathbf{X}^\top \mathbf{X}, \quad \alpha = E\{W_h(U)U^2\},$$

donde $U = Z^2$, $Z \sim \text{EC}(0, 1; h)$ con h la función generadora de densidad asociada a la mezcla de escala normal.

Apéndice A

Elementos de Álgebra Matricial

En este Apéndice se introduce la notación, definiciones y resultados básicos de álgebra lineal y matricial, esenciales para el estudio de modelos estadísticos multivariados y de regresión lineal. El material presentado a continuación puede ser hallado en textos como [Graybill \(1983\)](#), [Ravishanker y Dey \(2002\)](#) y [Magnus y Neudecker \(2007\)](#).

A.1. Vectores y matrices

Sea \mathbb{R}^n el espacio Euclidiano n -dimensional, de este modo $\mathbf{x} \in \mathbb{R}^n$ representa la n -upla

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix},$$

de números reales. Note que \mathbf{x} está *orientado* como un vector “columna”, y por tanto la transpuesta de \mathbf{x} es un vector fila,

$$\mathbf{x} = (x_1, \dots, x_n)^\top.$$

Una matriz $\mathbf{A} \in \mathbb{R}^{m \times n}$ es un arreglo de números reales

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix},$$

y escribimos $\mathbf{A} = (a_{ij})$. Los números reales a_{ij} son llamados elementos de \mathbf{A} .

A.2. Definiciones básicas y propiedades

La suma de dos matrices del mismo orden es definida como

$$\mathbf{A} + \mathbf{B} = (a_{ij}) + (b_{ij}) = (a_{ij} + b_{ij}),$$

el producto de una matriz por un escalar λ es

$$\lambda \mathbf{A} = \mathbf{A} \lambda = (\lambda a_{ij})$$

RESULTADO A.1 (Propiedades de la suma matricial). *Sean \mathbf{A}, \mathbf{B} y \mathbf{C} matrices del mismo orden y λ, μ escalares. Entonces:*

- (a) $\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$,
- (b) $(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$,
- (c) $(\lambda + \mu)\mathbf{A} = \lambda\mathbf{A} + \mu\mathbf{A}$,
- (d) $\lambda(\mathbf{A} + \mathbf{B}) = \lambda\mathbf{A} + \lambda\mathbf{B}$,
- (e) $\lambda\mu\mathbf{A} = (\lambda\mu)\mathbf{A}$.

Una matriz cuyos elementos son todos cero se denomina *matriz nula* y se denota por $\mathbf{0}$. Tenemos que

$$\mathbf{A} + (-1)\mathbf{A} = \mathbf{0}.$$

Si \mathbf{A} y \mathbf{B} son matrices $m \times n$ y $n \times p$, respectivamente, se define el producto de \mathbf{A} y \mathbf{B} como

$$\mathbf{AB} = \mathbf{C}, \quad \text{donde,} \quad c_{ij} = \sum_{k=1}^n a_{ik}b_{kj},$$

para $i = 1, \dots, m$ y $j = 1, \dots, p$.

RESULTADO A.2 (Propiedades del producto de matrices). Sean \mathbf{A}, \mathbf{B} y \mathbf{C} matrices de órdenes apropiados. Entonces:

- (a) $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$,
- (b) $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$,
- (c) $(\mathbf{A} + \mathbf{B})\mathbf{C} = \mathbf{AC} + \mathbf{BC}$.

Note que la existencia de \mathbf{AB} no implica la existencia de \mathbf{BA} y cuando ambos productos existen, en general no son iguales.

La transpuesta de una matriz $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{m \times n}$ es la matriz $n \times m$, \mathbf{A}^\top cuyo elemento ij está dado por a_{ji} , esto es

$$\mathbf{A}^\top = (a_{ji}).$$

RESULTADO A.3 (Propiedades de la transpuesta). Tenemos

- (a) $(\mathbf{A}^\top)^\top = \mathbf{A}$,
- (b) $(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$,
- (c) $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$.

Definimos el *producto interno* entre dos vectores $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ como

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i.$$

asociado al producto interno tenemos la norma Euclidian (o largo) de un vector \mathbf{x} definida como

$$\|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{1/2} = \left(\sum_{i=1}^n x_i^2 \right)^{1/2},$$

finalmente, la distancia Euclidian entre dos vectores \mathbf{a} y \mathbf{b} se define como

$$d(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\|.$$

RESULTADO A.4 (Propiedades del producto interno). Sean \mathbf{a}, \mathbf{b} y \mathbf{c} vectores n -dimensionales y λ un escalar, entonces

- (a) $\langle \mathbf{a}, \mathbf{b} \rangle = \langle \mathbf{b}, \mathbf{a} \rangle$,
- (b) $\langle \mathbf{a}, \mathbf{b} + \mathbf{c} \rangle = \langle \mathbf{a}, \mathbf{b} \rangle + \langle \mathbf{a}, \mathbf{c} \rangle$,
- (c) $\lambda \langle \mathbf{a}, \mathbf{b} \rangle = \langle \lambda \mathbf{a}, \mathbf{b} \rangle = \langle \mathbf{a}, \lambda \mathbf{b} \rangle$,
- (d) $\langle \mathbf{a}, \mathbf{a} \rangle \geq 0$ con la igualdad sólo si $\mathbf{a} = \mathbf{0}$,
- (e) $\|\mathbf{a} \pm \mathbf{b}\|^2 = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 \pm 2\langle \mathbf{a}, \mathbf{b} \rangle$,
- (f) $\|\mathbf{a} + \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$.

PROPOSICIÓN A.5 (Desigualdad de Cauchy-Schwarz). $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ con la igualdad sólo si $\mathbf{x} = \lambda \mathbf{y}$, para algún $\lambda \in \mathbb{R}$.

DEMOSTRACIÓN. Si $\mathbf{x} = \lambda \mathbf{y}$, el resultado es inmediato. Sino, note que

$$0 < \|\mathbf{x} - \lambda \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \lambda^2 \|\mathbf{y}\|^2 - 2\lambda \langle \mathbf{x}, \mathbf{y} \rangle, \quad \forall \lambda \in \mathbb{R},$$

de este modo el discriminante del polinomio cuadrático debe satisfacer $4\langle \mathbf{x}, \mathbf{y} \rangle^2 - 4\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 < 0$. \square

El ángulo θ entre dos vectores no nulos \mathbf{x}, \mathbf{y} se define en términos de su producto interno como

$$\cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{\mathbf{x}^\top \mathbf{y}}{\sqrt{\mathbf{x}^\top \mathbf{x}} \sqrt{\mathbf{y}^\top \mathbf{y}}},$$

dos vectores se dicen *ortogonales* sólo si $\mathbf{x}^\top \mathbf{y} = 0$.

El *producto externo* entre dos vectores $\mathbf{x} \in \mathbb{R}^m$ y $\mathbf{y} \in \mathbb{R}^n$ es la matriz $m \times n$

$$\mathbf{x} \wedge \mathbf{y} = \mathbf{x} \mathbf{y}^\top = (x_i y_j).$$

Una matriz se dice cuadrada si tiene el mismo número de filas que de columnas, una matriz cuadrada \mathbf{A} es triangular inferior (superior) si $a_{ij} = 0$ para $i < j$ (si $a_{ij} = 0$ para $i > j$). Una matriz cuadrada $\mathbf{A} = (a_{ij})$ se dice *simétrica* si $\mathbf{A}^\top = \mathbf{A}$ y *sesgo-simétrica* si $\mathbf{A}^\top = -\mathbf{A}$. Para cualquier matriz cuadrada $\mathbf{A} = (a_{ij})$ se define $\text{diag}(\mathbf{A})$ como

$$\text{diag}(\mathbf{A}) = \begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & a_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{pmatrix} = \text{diag}(a_{11}, a_{22}, \dots, a_{nn}).$$

Si $\mathbf{A} = \text{diag}(\mathbf{A})$, decimos que \mathbf{A} es *matriz diagonal*. Un tipo particular de matriz diagonal es la identidad

$$\mathbf{I} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} = (\delta_{ij}),$$

donde $\delta_{ij} = 1$ si $i = j$ y $\delta_{ij} = 0$ si $i \neq j$ (δ_{ij} se denomina *delta de Kronecker*). Tenemos que para $\mathbf{A} \in \mathbb{R}^{m \times n}$

$$\mathbf{I}_m \mathbf{A} = \mathbf{A} \mathbf{I}_n = \mathbf{A}.$$

Una matriz cuadrada se dice *ortogonal* si

$$\mathbf{A} \mathbf{A}^\top = \mathbf{A}^\top \mathbf{A} = \mathbf{I}$$

y sus columnas son ortonormales. Note que, si

$$\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n) \quad \text{con } \mathbf{a}_j \in \mathbb{R}^n,$$

entonces \mathbf{A} tiene columnas *ortonormales* si

$$\mathbf{a}_i^\top \mathbf{a}_j = \begin{cases} 1, & \text{si } i = j, \\ 0, & \text{si } i \neq j, \end{cases} \quad i, j = 1, \dots, n.$$

Una matriz rectangular $\mathbf{A} \in \mathbb{R}^{m \times n}$ puede tener la propiedad $\mathbf{A} \mathbf{A}^\top = \mathbf{I}_m$ ó $\mathbf{A}^\top \mathbf{A} = \mathbf{I}_n$ pero no ambas, en cuyo caso tal matriz se denomina semi-ortogonal.

Una matriz $\mathbf{A} \in \mathbb{R}^{n \times n}$, se dice idempotente si $\mathbf{A}^2 = \mathbf{A}$. Decimos que \mathbf{A} es *matriz de proyección* si es simétrica e idempotente, esto es, $\mathbf{A}^\top = \mathbf{A}$ y $\mathbf{A}^2 = \mathbf{A}$. Considere el siguiente resultado

RESULTADO A.6. *Suponga \mathbf{A} matriz $m \times m$, simétrica e idempotente. Entonces,*

- (a) $a_{ii} \geq 0$, $i = 1, \dots, n$.
- (b) $a_{ii} \leq 1$, $i = 1, \dots, n$.
- (c) $a_{ij} = a_{ji} = 0$, para todo $j \neq i$, si $a_{ii} = 0$ o $a_{ii} = 1$.

DEMOSTRACIÓN. Como \mathbf{A} es simétrica e idempotente, tenemos

$$\mathbf{A} = \mathbf{A}^2 = \mathbf{A}^\top \mathbf{A},$$

de ahí que

$$a_{ii} = \sum_{j=1}^n a_{ji}^2,$$

que claramente es no negativo. Además, podemos escribir

$$a_{ii} = a_{ii}^2 + \sum_{j \neq i} a_{ji}^2.$$

Por tanto, $a_{ii} \geq a_{ii}^2$ y de este modo (b) es satisfecha. Si $a_{ii} = 0$ o bien $a_{ii} = 1$, entonces $a_{ii} = a_{ii}^2$ y debemos tener

$$\sum_{j \neq i} a_{ji}^2 = 0,$$

lo que junto con la simetría de \mathbf{A} , establece (c). □

Cualquier matriz \mathbf{B} satisfaciendo

$$\mathbf{B}^2 = \mathbf{A}$$

se dice *raíz cuadrada* de \mathbf{A} y se denota como $\mathbf{A}^{1/2}$ tal matriz *no* necesita ser única.

A.2.1. Formas lineales y cuadráticas. Sea $\mathbf{a} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{n \times n}$ y $\mathbf{B} \in \mathbb{R}^{n \times m}$. La expresión $\mathbf{a}^\top \mathbf{x}$ se dice una *forma lineal* en \mathbf{x} y $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ una *forma cuadrática*, mientras que $\mathbf{x}^\top \mathbf{B} \mathbf{y}$ es una forma bilineal.

Sin pérdida de generalidad se asumirá que la matriz asociada a la forma cuadrática $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ es simétrica. Note que siempre es posible

$$\mathbf{x}^\top \mathbf{B} \mathbf{x} = \frac{1}{2} \mathbf{x}^\top (\mathbf{A}^\top + \mathbf{A}) \mathbf{x},$$

en cuyo caso tenemos que \mathbf{B} es matriz simétrica.

Decimos que una matriz simétrica \mathbf{A} es definida positiva (negativa) si $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$ ($\mathbf{x}^\top \mathbf{A} \mathbf{x} < 0$) para todo $\mathbf{x} \neq \mathbf{0}$. Cuando $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$ ($\mathbf{x}^\top \mathbf{A} \mathbf{x} \leq 0$) $\forall \mathbf{x}$ decimos que \mathbf{A} es semidefinida positiva (negativa).

Note que las matrices $\mathbf{B}^\top \mathbf{B}$ y $\mathbf{B} \mathbf{B}^\top$ son semidefinidas positivas y que \mathbf{A} es (semi)definida negativa sólo si $-\mathbf{A}$ es (semi)definida positiva.

RESULTADO A.7. *Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times p}$ y $\mathbf{C} \in \mathbb{R}^{n \times p}$ y \mathbf{x} vector n -dimensional. Entonces*

- (a) $\mathbf{A} \mathbf{x} = \mathbf{0} \Leftrightarrow \mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{0}$,
- (b) $\mathbf{A} \mathbf{B} = \mathbf{0} \Leftrightarrow \mathbf{A}^\top \mathbf{A} \mathbf{B} = \mathbf{0}$,
- (c) $\mathbf{A}^\top \mathbf{A} \mathbf{B} = \mathbf{A}^\top \mathbf{A} \mathbf{C} \Leftrightarrow \mathbf{A} \mathbf{B} = \mathbf{A} \mathbf{C}$.

DEMOSTRACIÓN. (a) Claramente $\mathbf{Ax} = \mathbf{0} \Rightarrow \mathbf{A}^\top \mathbf{Ax} = \mathbf{0}$. Por otro lado, si $\mathbf{A}^\top \mathbf{Ax} = \mathbf{0}$, entonces $\mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} = (\mathbf{Ax})^\top \mathbf{Ax} = 0$ y de ahí que $\mathbf{Ax} = \mathbf{0}$. (b) sigue desde (a). Finalmente, (c) sigue desde (b) mediante substituir $\mathbf{B} - \mathbf{C}$ por \mathbf{B} en (c). \square

RESULTADO A.8. Sean $\mathbf{A} \in \mathbb{R}^{m \times n}$ y \mathbf{B}, \mathbf{C} matrices $n \times n$ con \mathbf{B} simétrica. Entonces

- (a) $\mathbf{Ax} = \mathbf{0}, \forall \mathbf{x} \in \mathbb{R}^n$ sólo si $\mathbf{A} = \mathbf{0}$,
- (b) $\mathbf{x}^\top \mathbf{Bx} = 0, \forall \mathbf{x} \in \mathbb{R}^n$ sólo si $\mathbf{B} = \mathbf{0}$,
- (c) $\mathbf{x}^\top \mathbf{Cx} = 0, \forall \mathbf{x} \in \mathbb{R}^n$ sólo si $\mathbf{C}^\top = -\mathbf{C}$.

A.2.2. Rango de una matriz. Un conjunto de vectores $\mathbf{x}_1, \dots, \mathbf{x}_n$ se dice *linealmente independiente* si $\sum_i \alpha_i \mathbf{x}_i = \mathbf{0}$ implica que todos los $\alpha_i = 0$. Si $\mathbf{x}_1, \dots, \mathbf{x}_n$ no son linealmente independientes, ellos se dicen *linealmente dependientes*.

Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$, el *rango* columna (fila) de \mathbf{A} es el número de columnas (filas) linealmente independientes. Denotamos el rango de \mathbf{A} como

$$\text{rg}(\mathbf{A}),$$

note que

$$\text{rg}(\mathbf{A}) \leq \min(m, n).$$

Si $\text{rg}(\mathbf{A}) = n$ decimos que \mathbf{A} tiene rango columna completo. Si $\text{rg}(\mathbf{A}) = 0$, entonces \mathbf{A} es la matriz nula. Por otro lado, si $\mathbf{A} = \mathbf{0}$, entonces $\text{rg}(\mathbf{A}) = 0$.

RESULTADO A.9 (Propiedades del rango). Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$ y \mathbf{B}, \mathbf{C} matrices de órdenes apropiados, entonces

- (a) $\text{rg}(\mathbf{A}) = \text{rg}(\mathbf{A}^\top) = \text{rg}(\mathbf{A}^\top \mathbf{A}) = \text{rg}(\mathbf{AA}^\top)$,
- (b) $\text{rg}(\mathbf{AB}) \leq \min\{\text{rg}(\mathbf{A}), \text{rg}(\mathbf{B})\}$,
- (c) $\text{rg}(\mathbf{BAC}) = \text{rg}(\mathbf{A})$ si \mathbf{B} y \mathbf{C} son matrices de rango completo,
- (d) $\text{rg}(\mathbf{A} + \mathbf{B}) \leq \text{rg}(\mathbf{A}) + \text{rg}(\mathbf{B})$,
- (e) si $\mathbf{A} \in \mathbb{R}^{m \times n}$ y $\mathbf{Ax} = \mathbf{0}$ para algún $\mathbf{x} \neq \mathbf{0}$, entonces $\text{rg}(\mathbf{A}) \leq n - 1$.

El *espacio columna* de $\mathbf{A} \in \mathbb{R}^{m \times n}$, denotado por $\mathcal{M}(\mathbf{A})$, es el conjunto de vectores

$$\mathcal{M}(\mathbf{A}) = \{\mathbf{y} : \mathbf{y} = \mathbf{Ax} \text{ para algún } \mathbf{x} \in \mathbb{R}^n\}.$$

De este modo, $\mathcal{M}(\mathbf{A})$ es el espacio vectorial generado por las columnas de \mathbf{A} . La dimensión de este espacio es $\text{rg}(\mathbf{A})$. Se tiene que

$$\mathcal{M}(\mathbf{A}) = \mathcal{M}(\mathbf{AA}^\top)$$

para cualquier matriz \mathbf{A} .

El *espacio nulo*, $\mathcal{N}(\mathbf{A})$, de una matriz $\mathbf{A} \in \mathbb{R}^{m \times n}$ consiste de todos los vectores n -dimensionales \mathbf{x} , tal que $\mathbf{Ax} = \mathbf{0}$, esto es,

$$\mathcal{N}(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^n \text{ tal que } \mathbf{Ax} = \mathbf{0}\}.$$

Note que, el espacio nulo es el conjunto de todas las soluciones del sistema lineal homogéneo $\mathbf{Ax} = \mathbf{0}$. $\mathcal{N}(\mathbf{A})$ es un subespacio de \mathbb{R}^n y su dimensión se denomina *nulidad* de \mathbf{A} . Además $\mathcal{N}(\mathbf{A}) = \{\mathcal{M}(\mathbf{A})\}^\perp$. Finalmente, considere la siguiente proposición

RESULTADO A.10. Para cualquier matriz $\mathbf{A} \in \mathbb{R}^{m \times n}$, entonces $n = \dim(\mathcal{N}(\mathbf{A})) + \text{rg}(\mathbf{A})$.

A.2.3. Matriz inversa. Sea \mathbf{A} una matriz cuadrada de orden $n \times n$. Decimos que \mathbf{A} es *no singular* si $\text{rg}(\mathbf{A}) = n$, y que \mathbf{A} es *singular* si $\text{rg}(\mathbf{A}) < n$. De este modo, si \mathbf{A} es no singular, entonces existe una matriz no singular \mathbf{B} tal que

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I}_n.$$

La matriz \mathbf{B} , denotada \mathbf{A}^{-1} es única y se denomina *inversa* de \mathbf{A} .

RESULTADO A.11 (Propiedades de la inversa). *Siempre que todas las matrices inversas involucradas existan, tenemos que*

- (a) $(\mathbf{A}^{-1})^\top = (\mathbf{A}^\top)^{-1}$.
- (b) $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$.
- (c) $(\lambda\mathbf{A})^{-1} = \frac{1}{\lambda}\mathbf{A}^{-1}$.
- (d) $\mathbf{P}^{-1} = \mathbf{P}^\top$, si \mathbf{P} es matriz ortogonal.
- (e) Si $\mathbf{A} > 0$, entonces $\mathbf{A}^{-1} > 0$.
- (f) (Teorema de Sherman-Morrison-Woodbury)

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B})^{-1}\mathbf{DA}^{-1},$$

donde $\mathbf{A}, \mathbf{B}, \mathbf{C}$ y \mathbf{D} son matrices $m \times m$, $m \times n$, $n \times n$ y $n \times m$, respectivamente.

- (g) Si $1 \pm \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u} \neq 0$, entonces

$$(\mathbf{A} \pm \mathbf{uv}^\top)^{-1} = \mathbf{A}^{-1} \mp \frac{\mathbf{A}^{-1} \mathbf{uv}^\top \mathbf{A}^{-1}}{1 \pm \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u}},$$

es conocida como la fórmula de Sherman-Morrison.

- (h) $(\mathbf{I} + \lambda\mathbf{A})^{-1} = \mathbf{I} + \sum_{i=1}^{\infty} (-1)^i \lambda^i \mathbf{A}^i$.

A.2.4. Determinante de una matriz. El determinante de una matriz corresponde a la función $\det : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$, denotada comúnmente como $|\mathbf{A}| = \det(\mathbf{A})$ y definida como

$$|\mathbf{A}| = \sum (-1)^{\sigma(j_1, \dots, j_n)} \prod_{i=1}^n a_{ij_i}$$

donde la sumatoria es tomada sobre todas las permutaciones (j_1, \dots, j_n) del conjunto de enteros $(1, \dots, n)$, y $\sigma(j_1, \dots, j_n)$ es el número de transposiciones necesarias para cambiar $(1, \dots, n)$ en (j_1, \dots, j_n) (una transposición consiste en intercambiar dos números).

Una *submatriz* de \mathbf{A} es un arreglo rectangular obtenido mediante eliminar filas y columnas de \mathbf{A} . Un *menor* es el determinante de una submatriz cuadrada de \mathbf{A} . El menor asociado al elemento a_{ij} es el determinante de la submatriz de \mathbf{A} obtenida por eliminar su i -ésima fila y j -ésima columna. El *cofactor* de a_{ij} , digamos c_{ij} es $(-1)^{i+j}$ veces el menor de a_{ij} . La matriz $\mathbf{C} = (c_{ij})$ se denomina matriz cofactor de \mathbf{A} . La transpuesta de \mathbf{C} es llamada *adjunta* de \mathbf{A} y se denota $\mathbf{A}^\#$. Tenemos que

$$|\mathbf{A}| = \sum_{j=1}^n a_{ij} c_{ij} = \sum_{j=1}^n a_{jk} c_{jk}, \quad \text{para } i, k = 1, \dots, n.$$

RESULTADO A.12 (Propiedades del determinante). *Sea $\mathbf{A} \in \mathbb{R}^{n \times n}$ y λ un escalar. Entonces*

- (a) $|\mathbf{A}| = |\mathbf{A}^\top|$.
- (b) $|\mathbf{AB}| = |\mathbf{A}| |\mathbf{B}|$.

- (c) $|\lambda \mathbf{A}| = \lambda^n |\mathbf{A}|$.
- (d) $|\mathbf{A}^{-1}| = |\mathbf{A}|^{-1}$, si \mathbf{A} es no singular.
- (e) Si \mathbf{A} es matriz triangular, entonces $|\mathbf{A}| = \prod_{i=1}^n a_{ii}$.
- (f) El resultado en (e) también es válido para $\mathbf{A} = \text{diag}(\mathbf{A})$. Además, es evidente que $|\mathbf{I}_n| = 1$.
- (g) Si $\mathbf{A} \in \mathbb{R}^{m \times n}$ y $\mathbf{B} \in \mathbb{R}^{n \times m}$, entonces $|\mathbf{I}_m + \mathbf{AB}| = |\mathbf{I}_n + \mathbf{BA}|$.

A.2.5. La traza de una matriz. La traza de una matriz cuadrada $\mathbf{A} \in \mathbb{R}^{n \times n}$, denotada por $\text{tr}(\mathbf{A})$, es la suma de sus elementos diagonales:

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii}.$$

RESULTADO A.13 (Propiedades de la traza). *Siempre que las operaciones matriciales están definidas*

- (a) $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$,
- (b) $\text{tr}(\lambda \mathbf{A}) = \lambda \text{tr}(\mathbf{A})$ si λ es un escalar,
- (c) $\text{tr}(\mathbf{A}^\top) = \text{tr}(\mathbf{A})$,
- (d) $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ (propiedad cíclica de la traza),
- (e) $\text{tr}(\mathbf{A}) = 0$ si $\mathbf{A} = \mathbf{0}$.

Note en (d) que aunque ambas \mathbf{AB} y \mathbf{BA} son cuadradas, no necesitan ser del mismo orden.

Además, es directo que la normal vectorial (Euclidiana), satisface

$$\|\mathbf{x}\| = (\mathbf{x}^\top \mathbf{x})^{1/2} = (\text{tr} \mathbf{x} \mathbf{x}^\top)^{1/2},$$

de este modo, podemos definir una normal matricial (Euclidiana) como

$$\|\mathbf{A}\| = (\text{tr} \mathbf{A}^\top \mathbf{A})^{1/2}.$$

En efecto, se tiene que $\text{tr}(\mathbf{A}^\top \mathbf{A}) \geq 0$ con la igualdad sólo si $\mathbf{A} = \mathbf{0}$.

A.2.6. Valores y vectores propios. Si \mathbf{A} y \mathbf{B} son matrices reales del mismo orden, una matriz compleja \mathbf{Z} puede ser definida como

$$\mathbf{Z} = \mathbf{A} + i\mathbf{B},$$

donde i denota la unidad imaginaria que satisface $i^2 = -1$. El conjugado complejo de \mathbf{Z} , denotado por \mathbf{Z}^H , se define como

$$\mathbf{Z}^H = \mathbf{A}^\top - i\mathbf{B}^\top.$$

Una matriz $\mathbf{Z} \in \mathbb{C}^{n \times n}$ se dice *Hermitiana* si $\mathbf{Z}^H = \mathbf{Z}$ (equivalente complejo de una matriz simétrica) y *unitaria* si $\mathbf{Z}^H \mathbf{Z} = \mathbf{I}$ (equivalente complejo de una matriz ortogonal).

Sea \mathbf{A} una matriz cuadrada $n \times n$. Los *valores propios* de \mathbf{A} son definidos como las raíces de la *ecuación característica*

$$|\lambda \mathbf{I} - \mathbf{A}| = 0,$$

la ecuación anterior tiene n raíces, en general complejas y posiblemente con algunas repeticiones (multiplicidad). Sea λ un valor propio de \mathbf{A} , entonces existe un vector $\mathbf{v} \neq \mathbf{0} \in \mathbb{C}^n$ tal que $(\lambda \mathbf{I} - \mathbf{A})\mathbf{v} = \mathbf{0}$, esto es,

$$\mathbf{A}\mathbf{v} = \lambda \mathbf{v}.$$

el vector \mathbf{v} se denomina *vector propio* asociado al valor propio λ . Note que, si \mathbf{v} es un vector propio, también lo es $\alpha\mathbf{v}$, $\forall \alpha \in \mathbb{C}$, y en particular $\mathbf{v}/\|\mathbf{v}\|$ es un vector propio normalizado.

RESULTADO A.14. Si $\mathbf{A} \in \mathbb{C}^{n \times n}$ es matriz Hermitiana, entonces todos sus valores propios son reales

RESULTADO A.15. Si \mathbf{A} es matriz cuadrada $n \times n$ y \mathbf{G} es matriz no singular $n \times n$, entonces \mathbf{A} y $\mathbf{G}^{-1}\mathbf{A}\mathbf{G}$ tienen el mismo conjunto de valores propios (con las mismas multiplicidades)

DEMOSTRACIÓN. Note que

$$|\lambda\mathbf{I} - \mathbf{G}^{-1}\mathbf{A}\mathbf{G}| = |\lambda\mathbf{G}^{-1}\mathbf{G} - \mathbf{G}^{-1}\mathbf{A}\mathbf{G}| = |\mathbf{G}^{-1}||\lambda\mathbf{I} - \mathbf{A}||\mathbf{G}| = |\lambda\mathbf{I} - \mathbf{A}|$$

□

RESULTADO A.16. Una matriz singular tiene al menos un valor propio cero

DEMOSTRACIÓN. Si \mathbf{A} es matriz singular, entonces $\mathbf{A}\mathbf{v} = \mathbf{0}$ para algún $\mathbf{v} \neq \mathbf{0}$, luego desde $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$, tenemos que $\lambda = 0$. □

RESULTADO A.17. Una matriz simétrica es definida positiva (semidefinida positiva) sólo si todos sus valores propios son positivos (no-negativos).

DEMOSTRACIÓN. Si \mathbf{A} es definida positiva y $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$, entonces $\mathbf{v}^\top \mathbf{A}\mathbf{v} = \lambda\mathbf{v}^\top \mathbf{v}$. Ahora, como $\mathbf{v}^\top \mathbf{A}\mathbf{v} > 0$ y $\mathbf{v}^\top \mathbf{v} > 0$ implica $\lambda > 0$. La converso no será probada aquí. □

RESULTADO A.18. Una matriz idempotente sólo tiene valores propios 0 ó 1. Todos los valores propios de una matriz unitaria tienen modulo 1

DEMOSTRACIÓN. Sea \mathbf{A} matriz idempotente, esto es, $\mathbf{A}^2 = \mathbf{A}$. De este modo, si $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$, entonces

$$\lambda\mathbf{v} = \mathbf{A}\mathbf{v} = \mathbf{A}^2\mathbf{v} = \lambda\mathbf{A}\mathbf{v} = \lambda^2\mathbf{v}$$

y de ahí que $\lambda = \lambda^2$, esto implica que $\lambda = 0$ ó $\lambda = 1$.

Por otro lado, si \mathbf{A} es unitaria, entonces $\mathbf{A}^H\mathbf{A} = \mathbf{I}$. De este modo, si $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$, entonces

$$\mathbf{v}^H\mathbf{A}^H = \bar{\lambda}\mathbf{v}^H,$$

luego

$$\mathbf{v}^H\mathbf{v} = \mathbf{v}^H\mathbf{A}^H\mathbf{A}\mathbf{v} = \bar{\lambda}\lambda\mathbf{v}^H\mathbf{v}.$$

Como $\mathbf{v}^H\mathbf{v} \neq 0$, obtenemos que $\bar{\lambda}\lambda = 1$ y de ahí que $|\lambda| = 1$. □

RESULTADO A.19 (Propiedades de la matrices idempotentes). Sea \mathbf{A} matriz $n \times n$, entonces

- (a) \mathbf{A}^\top y $\mathbf{I} - \mathbf{A}$ son idempotentes sólo si \mathbf{A} es idempotente,
- (b) si \mathbf{A} es idempotente, entonces $\text{rg}(\mathbf{A}) = \text{tr}(\mathbf{A}) = r$. Si $\text{rg}(\mathbf{A}) = n$, entonces $\mathbf{A} = \mathbf{I}$.

RESULTADO A.20. Si $\mathbf{A} \in \mathbb{C}^{n \times n}$ es matriz Hermitiana y $\mathbf{v}_1, \mathbf{v}_2$ son vectores propios asociados a λ_1 y λ_2 , respectivamente, donde $\lambda_1 \neq \lambda_2$. Entonces $\mathbf{v}_1 \perp \mathbf{v}_2$.

El resultado anterior muestra que si todos los valores propios de una matriz Hermitiana \mathbf{A} son distintos, entonces existe una base ortonormal de vectores propios tal que \mathbf{A} es diagonalizable.

PROPOSICIÓN A.21 (Descomposición de Schur). Sea $\mathbf{A} \in \mathbb{C}^{n \times n}$. Entonces existe una matriz unitaria $\mathbf{U} \in \mathbb{C}^{n \times n}$ y una matriz triangular \mathbf{M} cuyos elementos diagonales son los valores propios de \mathbf{A} , tal que

$$\mathbf{U}^H \mathbf{A} \mathbf{U} = \mathbf{M}.$$

PROPOSICIÓN A.22 (Descomposición espectral). Sea $\mathbf{A} \in \mathbb{C}^{n \times n}$ matriz Hermitiana. Entonces existe una matriz unitaria $\mathbf{U} \in \mathbb{C}^{n \times n}$ tal que

$$\mathbf{U}^H \mathbf{A} \mathbf{U} = \mathbf{\Lambda},$$

donde $\mathbf{\Lambda} = \text{diag}(\boldsymbol{\lambda})$ es matriz diagonal cuyos elementos diagonales son los valores propios de \mathbf{A} .

Para aplicaciones en Estadística siempre haremos uso de la Proposición A.22 considerando \mathbf{A} matriz simétrica, en cuyo caso todos sus valores propios serán reales y \mathbf{U} será una matriz ortogonal. Para $\mathbf{Q} \in \mathbb{R}^{n \times n}$ matriz ortogonal, denotamos el grupo de matrices ortogonales como

$$\mathcal{O}_n = \{\mathbf{Q} \in \mathbb{R}^{n \times n} : \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}\}$$

Note que si \mathbf{A} es matriz simétrica y definida positiva, entonces

$$\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top = (\mathbf{U} \mathbf{\Lambda}^{1/2})(\mathbf{U} \mathbf{\Lambda}^{1/2})^\top = (\mathbf{U} \mathbf{\Lambda}^{1/2} \mathbf{U}^\top)^2$$

donde $\mathbf{\Lambda} = \text{diag}(\boldsymbol{\lambda})$ y $\mathbf{\Lambda}^{1/2} = \text{diag}(\boldsymbol{\lambda}^{1/2})$. Por tanto,

$$\mathbf{A} = \mathbf{M} \mathbf{M}^\top, \quad \text{con} \quad \mathbf{M} = \mathbf{U} \mathbf{\Lambda}^{1/2},$$

o bien,

$$\mathbf{A} = \mathbf{B}^2, \quad \text{con} \quad \mathbf{B} = \mathbf{U} \mathbf{\Lambda}^{1/2} \mathbf{U}^\top,$$

esto es, \mathbf{B} es una matriz raíz cuadrada de \mathbf{A} .

RESULTADO A.23. Sea \mathbf{A} matriz simétrica $n \times n$, con valores propios $\lambda_1, \dots, \lambda_n$. Entonces

- (a) $\text{tr}(\mathbf{A}) = \sum_{i=1}^n \lambda_i$,
- (b) $|\mathbf{A}| = \prod_{i=1}^n \lambda_i$.

DEMOSTRACIÓN. Usando que $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$. Tenemos

$$\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top) = \text{tr}(\mathbf{\Lambda} \mathbf{U}^\top \mathbf{U}) = \text{tr}(\mathbf{\Lambda}) = \sum_{i=1}^n \lambda_i$$

y

$$|\mathbf{A}| = |\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top| = |\mathbf{U}| |\mathbf{\Lambda}| |\mathbf{U}^\top| = |\mathbf{\Lambda}| = \prod_{i=1}^n \lambda_i$$

□

RESULTADO A.24. Si \mathbf{A} es una matriz simétrica con r valores propios distintos de cero, entonces $\text{rg}(\mathbf{A}) = r$.

DEMOSTRACIÓN. Tenemos que $\mathbf{U}^\top \mathbf{A} \mathbf{U} = \mathbf{\Lambda}$ y de ahí que

$$\text{rg}(\mathbf{A}) = \text{rg}(\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top) = \text{rg}(\mathbf{\Lambda}) = r$$

□

A.2.7. Matrices (semi)definidas positivas.

PROPOSICIÓN A.25. Sea \mathbf{A} matriz definida positiva y \mathbf{B} semidefinida positiva. Entonces

$$|\mathbf{A} + \mathbf{B}| \geq |\mathbf{A}|,$$

con la igualdad sólo si $\mathbf{B} = \mathbf{0}$.

DEMOSTRACIÓN. Tenemos $\mathbf{U}^\top \mathbf{A} \mathbf{U} = \mathbf{\Lambda}$, con $\mathbf{\Lambda} = \text{diag}(\boldsymbol{\lambda})$ y $\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{I}$. Luego,

$$\mathbf{A} + \mathbf{B} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top + \mathbf{B} = \mathbf{U} \mathbf{\Lambda}^{1/2} (\mathbf{I} + \mathbf{\Lambda}^{-1/2} \mathbf{U}^\top \mathbf{B} \mathbf{U} \mathbf{\Lambda}^{-1/2}) \mathbf{\Lambda}^{1/2} \mathbf{U}^\top,$$

de este modo

$$\begin{aligned} |\mathbf{A} + \mathbf{B}| &= |\mathbf{U} \mathbf{\Lambda}^{1/2}| |\mathbf{I} + \mathbf{\Lambda}^{-1/2} \mathbf{U}^\top \mathbf{B} \mathbf{U} \mathbf{\Lambda}^{-1/2}| |\mathbf{\Lambda}^{1/2} \mathbf{U}^\top| \\ &= |\mathbf{U} \mathbf{\Lambda}^{1/2} \mathbf{\Lambda}^{1/2} \mathbf{U}^\top| |\mathbf{I} + \mathbf{\Lambda}^{-1/2} \mathbf{U}^\top \mathbf{B} \mathbf{U} \mathbf{\Lambda}^{-1/2}| \\ &= |\mathbf{A}| |\mathbf{I} + \mathbf{\Lambda}^{-1/2} \mathbf{U}^\top \mathbf{B} \mathbf{U} \mathbf{\Lambda}^{-1/2}|. \end{aligned}$$

Si $\mathbf{B} = \mathbf{0}$, tenemos $|\mathbf{A} + \mathbf{B}| = |\mathbf{A}|$. Por otro lado, si $\mathbf{B} \neq \mathbf{0}$. Entonces la matriz $\mathbf{I} + \mathbf{\Lambda}^{-1/2} \mathbf{U}^\top \mathbf{B} \mathbf{U} \mathbf{\Lambda}^{-1/2}$ tendrá al menos un valor propio no nulo y por tanto, $|\mathbf{I} + \mathbf{\Lambda}^{-1/2} \mathbf{U}^\top \mathbf{B} \mathbf{U} \mathbf{\Lambda}^{-1/2}| > 1$, esto es $|\mathbf{A} + \mathbf{B}| > |\mathbf{A}|$. \square

Para dos matrices simétricas \mathbf{A} y \mathbf{B} , escribimos $\mathbf{A} \geq \mathbf{B}$ si $\mathbf{A} - \mathbf{B}$ es semidefinida positiva. Análogamente, escribimos $\mathbf{A} > \mathbf{B}$ si $\mathbf{A} - \mathbf{B}$ es definida positiva.

RESULTADO A.26. Sean \mathbf{A} , \mathbf{B} matrices definidas positivas $n \times n$. Entonces $\mathbf{A} > \mathbf{B}$ sólo si $\mathbf{B}^{-1} > \mathbf{A}^{-1}$.

PROPOSICIÓN A.27. Sean \mathbf{A} y \mathbf{B} matrices definidas positivas y $\mathbf{A} - \mathbf{B} \geq \mathbf{0}$. Entonces $|\mathbf{A}| \geq |\mathbf{B}|$ con la igualdad sólo si $\mathbf{A} = \mathbf{B}$.

DEMOSTRACIÓN. Sea $\mathbf{C} = \mathbf{A} - \mathbf{B}$. Como \mathbf{B} es definida positiva y \mathbf{C} es semidefinida positiva, tenemos por la Proposición A.25 que $|\mathbf{B} + \mathbf{C}| \geq |\mathbf{B}|$, con la igualdad sólo si $\mathbf{C} = \mathbf{0}$. \square

A.2.8. Descomposiciones matriciales.

PROPOSICIÓN A.28 (Descomposición LDL). Si $\mathbf{A} \in \mathbb{R}^{n \times n}$ es matriz simétrica y no singular, entonces existe \mathbf{L} matriz triangular inferior y $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$, tal que

$$\mathbf{A} = \mathbf{L} \mathbf{D} \mathbf{L}^\top.$$

PROPOSICIÓN A.29 (Descomposición Cholesky). Si $\mathbf{A} \in \mathbb{R}^{n \times n}$ es simétrica y definida positiva, entonces existe una única matriz triangular inferior $\mathbf{G} \in \mathbb{R}^{n \times n}$ (factor Cholesky) con elementos diagonales positivos, tal que

$$\mathbf{A} = \mathbf{G} \mathbf{G}^\top.$$

PROPOSICIÓN A.30 (Descomposición ortogonal-triangular). Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$, entonces existe $\mathbf{Q} \in \mathcal{O}_m$ y $\mathbf{R} \in \mathbb{R}^{m \times n}$, tal que

$$\mathbf{A} = \mathbf{Q} \mathbf{R},$$

donde

$$\mathbf{R} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix}$$

con $\mathbf{R}_1 \in \mathbb{R}^{n \times n}$ matriz triangular superior, aquí suponemos que $m \geq n$. Si $\text{rg}(\mathbf{A}) = r$, entonces las primeras n columnas de \mathbf{Q} forman una base ortonormal para $\mathcal{M}(\mathbf{A})$.

Note que, si $\mathbf{A} = \mathbf{Q}\mathbf{R}$ entonces

$$\mathbf{A}^\top \mathbf{A} = \mathbf{R}^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{R} = \mathbf{R}^\top \mathbf{R} = \mathbf{R}_1^\top \mathbf{R}_1,$$

y \mathbf{R}_1 corresponde al factor Cholesky de $\mathbf{A}^\top \mathbf{A}$.

PROPOSICIÓN A.31 (Descomposición valor singular). Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$ con $\text{rg}(\mathbf{A}) = r$, entonces existen matrices $\mathbf{U} \in \mathcal{O}_m$, $\mathbf{V} \in \mathcal{O}_n$, tal que

$$\mathbf{A} = \mathbf{U} \begin{pmatrix} \mathbf{D}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V}^\top,$$

donde $\mathbf{D}_r = \text{diag}(\delta_1, \dots, \delta_r)$ con $\delta_i > 0$ para $i = 1, \dots, r$, llamados **valores singulares** de \mathbf{A} .

A.2.9. Matrices particionadas. Sea \mathbf{A} una matriz $m \times n$. Considere particionar \mathbf{A} como sigue

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}, \quad (\text{A.1})$$

donde $\mathbf{A}_{11} \in \mathbb{R}^{m_1 \times n_1}$, $\mathbf{A}_{12} \in \mathbb{R}^{m_1 \times n_2}$, $\mathbf{A}_{21} \in \mathbb{R}^{m_2 \times n_1}$, $\mathbf{A}_{22} \in \mathbb{R}^{m_2 \times n_2}$, y $m_1 + m_2 = m$, $n_1 + n_2 = n$.

Sea $\mathbf{B} \in \mathbb{R}^{m \times n}$ particionada de manera análoga a \mathbf{A} , entonces

$$\mathbf{A} + \mathbf{B} = \begin{pmatrix} \mathbf{A}_{11} + \mathbf{B}_{11} & \mathbf{A}_{12} + \mathbf{B}_{12} \\ \mathbf{A}_{21} + \mathbf{B}_{21} & \mathbf{A}_{22} + \mathbf{B}_{22} \end{pmatrix}.$$

Ahora, considere $\mathbf{C} \in \mathbb{R}^{n \times p}$ particionada en submatrices \mathbf{C}_{ij} , para $i, j = 1, 2$ con dimensiones adecuadas, entonces

$$\mathbf{A}\mathbf{C} = \begin{pmatrix} \mathbf{A}_{11}\mathbf{C}_{11} + \mathbf{A}_{12}\mathbf{C}_{21} & \mathbf{A}_{11}\mathbf{C}_{12} + \mathbf{A}_{12}\mathbf{C}_{22} \\ \mathbf{A}_{21}\mathbf{C}_{11} + \mathbf{A}_{22}\mathbf{C}_{21} & \mathbf{A}_{21}\mathbf{C}_{12} + \mathbf{A}_{22}\mathbf{C}_{22} \end{pmatrix}.$$

La transpuesta de \mathbf{A} está dada por

$$\mathbf{A}^\top = \begin{pmatrix} \mathbf{A}_{11}^\top & \mathbf{A}_{21}^\top \\ \mathbf{A}_{12}^\top & \mathbf{A}_{22}^\top \end{pmatrix}.$$

Si \mathbf{A}_{12} y \mathbf{A}_{21} son matrices nulas y si ambas \mathbf{A}_{11} y \mathbf{A}_{22} son matrices no singulares, entonces la inversa de \mathbf{A} es

$$\mathbf{A}^{-1} = \begin{pmatrix} \mathbf{A}_{11}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{22}^{-1} \end{pmatrix}.$$

En general, si \mathbf{A} es matriz no singular particionada como en (A.1) y $\mathbf{D} = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$ también es no singular, entonces

$$\mathbf{A}^{-1} = \begin{pmatrix} \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{D}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{D}^{-1} \end{pmatrix}.$$

Por otro lado, si \mathbf{A} es no singular y $\mathbf{E} = \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$ es no singular, entonces

$$\mathbf{A}^{-1} = \begin{pmatrix} \mathbf{E}^{-1} & -\mathbf{E}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ -\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{E}^{-1} & \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{E}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \end{pmatrix}.$$

Considere el determinante

$$\begin{vmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{A}_{22} \end{vmatrix} = |\mathbf{A}_{11}| |\mathbf{A}_{22}| = \begin{vmatrix} \mathbf{A}_{11} & \mathbf{0} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{vmatrix},$$

si \mathbf{A}_{11} y \mathbf{A}_{22} son matrices cuadradas.

Ahora, para una matriz particionada como en (A.1) con $m_1 = n_1$ y $m_2 = n_2$, tenemos

$$|\mathbf{A}| = |\mathbf{A}_{11}| |\mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12}| = |\mathbf{A}_{22}| |\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21}|,$$

si \mathbf{A}_{11} y \mathbf{A}_{22} son matrices no singulares.

A.3. Inversa generalizada y sistemas de ecuaciones lineales

En esta sección se generaliza el concepto de invertibilidad para matrices singulares así como para matrices rectangulares. En particular, introducimos la inversa Moore-Penrose (MP), generalización que permite resolver de forma explícita un sistema de ecuaciones lineales.

A.3.1. Inversa Moore-Penrose. Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$, la inversa Moore-Penrose, $\mathbf{G} \in \mathbb{R}^{n \times m}$ debe satisfacer las siguientes condiciones

$$\mathbf{AGA} = \mathbf{A}, \quad (\text{A.2})$$

$$\mathbf{GAG} = \mathbf{G}, \quad (\text{A.3})$$

$$(\mathbf{AG})^\top = \mathbf{AG}, \quad (\text{A.4})$$

$$(\mathbf{GA})^\top = \mathbf{GA}. \quad (\text{A.5})$$

La inversa MP de \mathbf{A} se denota comunmente como \mathbf{A}^+ . Si \mathbf{G} satisface sólo la condición en (A.2) entonces decimos que \mathbf{G} es una inversa generalizada y la denotamos por \mathbf{A}^- .

PROPOSICIÓN A.32 (Unicidad de la inversa MP). *Para cada \mathbf{A} , existe una única \mathbf{A}^+ .*

RESULTADO A.33 (Propiedades de la inversa MP).

- (a) $\mathbf{A}^+ = \mathbf{A}^{-1}$ para \mathbf{A} matriz no singular,
- (b) $(\mathbf{A}^+)^+ = \mathbf{A}$,
- (c) $(\mathbf{A}^\top)^+ = (\mathbf{A}^+)^\top$,
- (d) $\mathbf{A}^+ = \mathbf{A}$ si \mathbf{A} es simétrica e idempotente,
- (e) \mathbf{AA}^+ y $\mathbf{A}^+\mathbf{A}$ son idempotentes,
- (f) $\text{rg}(\mathbf{A}) = \text{rg}(\mathbf{A}^+) = \text{rg}(\mathbf{AA}^+) = \text{rg}(\mathbf{A}^+\mathbf{A})$,
- (g) $\mathbf{A}^\top \mathbf{AA}^+ = \mathbf{A} = \mathbf{A}^+ \mathbf{AA}^\top$,
- (h) $\mathbf{A}^\top \mathbf{A}^{+\top} \mathbf{A}^+ = \mathbf{A}^+ = \mathbf{A}^+ \mathbf{A}^{+\top} \mathbf{A}^\top$,
- (i) $\mathbf{A}^+ = (\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top = \mathbf{A}^\top (\mathbf{AA}^\top)^+$,
- (j) $\mathbf{A}^+ = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$, si \mathbf{A} tiene rango columna completo,
- (k) $\mathbf{A}^+ = \mathbf{A}^\top (\mathbf{AA}^\top)^{-1}$, si \mathbf{A} tiene rango fila completo.

A.3.2. Solución de sistemas de ecuaciones lineales. La solución general de un sistema de ecuaciones homogéneo $\mathbf{Ax} = \mathbf{0}$ es

$$\mathbf{x} = (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \mathbf{q},$$

con \mathbf{q} un vector arbitrario. La solución de $\mathbf{Ax} = \mathbf{0}$ es única sólo si \mathbf{A} tiene rango columna completo, esto es, $\mathbf{A}^\top \mathbf{A}$ es no singular. El sistema homogéneo $\mathbf{Ax} = \mathbf{0}$ siempre tiene al menos una solución, digamos $\mathbf{x} = \mathbf{0}$.

El sistema no homogéneo

$$\mathbf{Ax} = \mathbf{b},$$

tendrá al menos una solución si es *consistente*.

PROPOSICIÓN A.34. Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$ y \mathbf{b} vector $m \times 1$. Entonces son equivalentes:

- (a) la ecuación $\mathbf{Ax} = \mathbf{b}$ tiene una solución para \mathbf{x} ,
- (b) $\mathbf{b} \in \mathcal{M}(\mathbf{A})$,
- (c) $\text{rg}(\mathbf{A} : \mathbf{b}) = \text{rg}(\mathbf{A})$,
- (d) $\mathbf{AA}^+\mathbf{b} = \mathbf{b}$.

PROPOSICIÓN A.35. Una condición necesaria y suficiente para que la ecuación $\mathbf{Ax} = \mathbf{b}$ tenga una solución es que

$$\mathbf{AA}^+\mathbf{b} = \mathbf{b},$$

en cuyo caso la solución general está dada por

$$\mathbf{x} = \mathbf{A}^+\mathbf{b} + (\mathbf{I} - \mathbf{A}^+\mathbf{A})\mathbf{q},$$

donde \mathbf{q} es un vector arbitrario.

Si el sistema $\mathbf{Ax} = \mathbf{b}$ es consistente, entonces tendrá solución única sólo si \mathbf{A} es de rango completo, en cuyo caso la solución está dada por $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$.

PROPOSICIÓN A.36. Una condición necesaria y suficiente para que la ecuación matricial $\mathbf{AXB} = \mathbf{C}$ tenga una solución es que

$$\mathbf{AA}^+\mathbf{CB}^+\mathbf{B} = \mathbf{C},$$

en cuyo caso la solución general es

$$\mathbf{X} = \mathbf{A}^+\mathbf{CB}^+ + \mathbf{Q} - \mathbf{A}^+\mathbf{AQB}^+\mathbf{B},$$

donde \mathbf{Q} es una matriz arbitraria de órdenes apropiados.

Apéndice B

Diferenciación matricial

En esta sección haremos uso de la siguiente notación. ϕ , \mathbf{f} y \mathbf{F} representan funciones escalar, vectorial y matricial, respectivamente mientras que ζ , \mathbf{x} y \mathbf{X} argumentos escalar, vectorial y matricial, respectivamente.

A partir de esta convención es directo que podemos escribir los siguientes casos particulares:

$$\begin{aligned}\phi(\zeta) &= \zeta^2, & \phi(\mathbf{x}) &= \mathbf{a}^\top \mathbf{x}, & \phi(\mathbf{X}) &= \text{tr}(\mathbf{X}^\top \mathbf{X}), \\ \mathbf{f}(\zeta) &= (\zeta, \zeta^2)^\top, & \mathbf{f}(\mathbf{x}) &= \mathbf{A}\mathbf{x}, & \mathbf{f}(\mathbf{X}) &= \mathbf{X}\mathbf{a}, \\ \mathbf{F}(\zeta) &= \zeta^2 \mathbf{I}_n, & \mathbf{F}(\mathbf{x}) &= \mathbf{x}\mathbf{x}^\top, & \mathbf{F}(\mathbf{X}) &= \mathbf{X}^\top.\end{aligned}$$

Existen varias definiciones para la derivada de una función matricial $\mathbf{F}(\mathbf{X})$ con relación a su argumento (matricial) \mathbf{X} . En este apéndice nos enfocamos en el cálculo diferencial propuesto por [Magnus y Neudecker \(1985\)](#).

Considere $\phi : S \rightarrow \mathbb{R}$ con $S \subset \mathbb{R}^n$, se define la derivada de ϕ con relación a $\mathbf{x} \in S$ como

$$\frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} = \left(\frac{\partial \phi}{\partial x_1}, \dots, \frac{\partial \phi}{\partial x_n} \right)^\top = \left(\frac{\partial \phi}{\partial x_i} \right) \in \mathbb{R}^n$$

de este modo, introducimos la notación

$$\mathbf{D}\phi(\mathbf{x}) = \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}^\top} \in \mathbb{R}^{1 \times n}.$$

Ahora, si $\mathbf{f} : S \rightarrow \mathbb{R}^m$, $S \subset \mathbb{R}^n$. Entonces la matriz $m \times n$,

$$\mathbf{D}\mathbf{f}(\mathbf{x}) = \begin{pmatrix} \mathbf{D}f_1(\mathbf{x}) \\ \vdots \\ \mathbf{D}f_m(\mathbf{x}) \end{pmatrix} = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}^\top},$$

es la *derivada* o *matriz Jacobiana* de \mathbf{f} . La transpuesta de la matriz Jacobiana $\mathbf{D}\mathbf{f}(\mathbf{x})$ se denomina *gradiente* de $\mathbf{f}(\mathbf{x})$.

B.1. Aproximación de primer orden

Considere la fórmula de Taylor de primer orden,

$$\phi(c + u) = \phi(c) + u\phi'(c) + r_c(u),$$

donde el *resto*

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

es de orden más pequeño que u conforme $u \rightarrow 0$. Note también que

$$\lim_{u \rightarrow 0} \frac{\phi(c + u) - \phi(c)}{u} = \phi'(c).$$

De este modo, se define

$$d\phi(c; u) = u\phi'(c),$$

como el (*primer*) *diferencial* de ϕ en c con incremento u . Esto motiva la siguiente definición.

DEFINICIÓN B.1 (Diferencial de una función vectorial). Sea $\mathbf{f} : S \rightarrow \mathbb{R}^m$, $S \subset \mathbb{R}^n$, si existe una matriz $\mathbf{A} \in \mathbb{R}^{m \times n}$, tal que

$$\mathbf{f}(\mathbf{c} + \mathbf{u}) = \mathbf{f}(\mathbf{c}) + \mathbf{A}(\mathbf{c})\mathbf{u} + \mathbf{r}_c(\mathbf{u}),$$

para todo $\mathbf{u} \in \mathbb{R}^n$ con $\|\mathbf{u}\| < \delta$, y

$$\lim_{\mathbf{u} \rightarrow 0} \frac{\mathbf{r}_c(\mathbf{u})}{\|\mathbf{u}\|} = \mathbf{0},$$

entonces la función \mathbf{f} se dice diferenciable en \mathbf{c} . El vector $m \times 1$

$$d\mathbf{f}(\mathbf{c}; \mathbf{u}) = \mathbf{A}(\mathbf{c})\mathbf{u},$$

se denomina primer diferencial de \mathbf{f} en \mathbf{c} con incremento \mathbf{u} .

[Magnus y Neudecker \(1985\)](#) mostraron la existencia y unicidad del diferencial $d\mathbf{f}(\mathbf{c}; \mathbf{u})$ de una función $\mathbf{f} : S \rightarrow \mathbb{R}^m$, $S \subset \mathbb{R}^n$ ($\mathbf{c} \in S$), dado por

$$d\mathbf{f}(\mathbf{c}; \mathbf{u}) = \mathbf{A}(\mathbf{c})\mathbf{u}$$

también mostraron la regla de la cadena e invarianza de Cauchy para el diferencial y enunciaron su primer teorema de identificación.

TEOREMA B.2 (Primer teorema de identificación). Sea $\mathbf{f} : S \rightarrow \mathbb{R}^m$, $S \subset \mathbb{R}^n$ función diferenciable, $\mathbf{c} \in S$ y \mathbf{u} un vector n -dimensional. Entonces

$$d\mathbf{f}(\mathbf{c}; \mathbf{u}) = (D\mathbf{f}(\mathbf{c}))\mathbf{u}.$$

La matriz $D\mathbf{f}(\mathbf{c}) \in \mathbb{R}^{m \times n}$ se denomina matriz Jacobiana. Tenemos también que

$$\nabla \mathbf{f}(\mathbf{c}) = (D\mathbf{f}(\mathbf{c}))^\top$$

es la matriz gradiente de \mathbf{f} .

Sea $\mathbf{f} : S \rightarrow \mathbb{R}^m$, $S \subset \mathbb{R}^n$ y $f_i : S \rightarrow \mathbb{R}$ el i -ésimo componente de \mathbf{f} ($i = 1, \dots, m$). Sea \mathbf{e}_j un vector n -dimensional cuyo j -ésimo elemento es uno y los restantes son cero, y considere

$$\lim_{t \rightarrow 0} \frac{f_i(\mathbf{c} + t\mathbf{e}_j) - f_i(\mathbf{c})}{t}$$

si el límite existe, se denomina la j -ésima *derivada parcial* de f_i en \mathbf{c} y es denotada por $D_j f_i(\mathbf{c})$. Note que el elemento ij de $D\mathbf{f}(\mathbf{c})$ es $D_j f_i(\mathbf{c})$.

B.2. Funciones matriciales

Considere algunos ejemplos de funciones matriciales

$$\mathbf{F}(\zeta) = \begin{pmatrix} \cos(\zeta) & \sin(\zeta) \\ -\sin(\zeta) & \cos(\zeta) \end{pmatrix}, \quad \mathbf{F}(\mathbf{x}) = \mathbf{x}\mathbf{x}^\top, \quad \mathbf{F}(\mathbf{X}) = \mathbf{X}^\top, \quad \mathbf{X} \in \mathbb{R}^{n \times q}.$$

Antes de considerar el diferencial de una función matricial $\mathbf{F} : S \rightarrow \mathbb{R}^{m \times p}$, $S \subset \mathbb{R}^{n \times q}$ introducimos dos conceptos preliminares: la vectorización de una matriz y el producto Kronecker.

DEFINICIÓN B.3 (Operador de vectorización). Sea $\mathbf{A} \in \mathbb{R}^{n \times q}$ particionada como

$$\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_q),$$

donde $\mathbf{a}_k \in \mathbb{R}^n$ es la k -ésima columna de \mathbf{A} . Entonces

$$\text{vec}(\mathbf{A}) = \begin{pmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_q \end{pmatrix}.$$

DEFINICIÓN B.4 (Producto Kronecker). Sea $\mathbf{A} \in \mathbb{R}^{m \times n}$ y $\mathbf{B} \in \mathbb{R}^{p \times q}$, entonces el producto Kronecker entre \mathbf{A} y \mathbf{B} denotado por $\mathbf{A} \otimes \mathbf{B}$ es la matriz $mp \times nq$ definida como

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & \dots & a_{1n}\mathbf{B} \\ \vdots & & \vdots \\ a_{m1}\mathbf{B} & \dots & a_{mn}\mathbf{B} \end{pmatrix}$$

RESULTADO B.5. Sean $\mathbf{A}, \mathbf{B}, \mathbf{C}$ y \mathbf{D} matrices de órdenes apropiados y λ escalar. Entonces

- (a) $\mathbf{A} \otimes \mathbf{B} \otimes \mathbf{C} = (\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}),$
- (b) $(\mathbf{A} + \mathbf{B}) \otimes (\mathbf{C} + \mathbf{D}) = \mathbf{A} \otimes \mathbf{C} + \mathbf{B} \otimes \mathbf{C} + \mathbf{A} \otimes \mathbf{D} + \mathbf{B} \otimes \mathbf{D},$
- (c) $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD},$
- (d) $\lambda \otimes \mathbf{A} = \lambda \mathbf{A} = \mathbf{A} \otimes \lambda,$
- (e) $(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top,$
- (f) $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1},$
- (g) $(\mathbf{A} \otimes \mathbf{B})^- = \mathbf{A}^- \otimes \mathbf{B}^-.$

RESULTADO B.6. Sean $\mathbf{A} \in \mathbb{R}^{n \times n}$ y $\mathbf{B} \in \mathbb{R}^{p \times p}$. Entonces

- (a) $\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A}) \text{tr}(\mathbf{B}),$
- (b) $|\mathbf{A} \otimes \mathbf{B}| = |\mathbf{A}|^p |\mathbf{B}|^n,$
- (c) $\text{rg}(\mathbf{A} \otimes \mathbf{B}) = \text{rg}(\mathbf{A}) \text{rg}(\mathbf{B}).$

Observe que, si $\mathbf{a} \in \mathbb{R}^n$ y $\mathbf{b} \in \mathbb{R}^p$, entonces

$$\mathbf{ab}^\top = \mathbf{a} \otimes \mathbf{b}^\top = \mathbf{b}^\top \otimes \mathbf{a},$$

por otro lado, tenemos que

$$\text{vec}(\mathbf{ab}^\top) = \text{vec}(\mathbf{a} \otimes \mathbf{b}^\top) = \text{vec}(\mathbf{b}^\top \otimes \mathbf{a}) = \mathbf{b} \otimes \mathbf{a}.$$

Estos resultados sugieren una conexión entre el operador de vectorización, el producto Kronecker y la traza. Considere el siguiente resultado

RESULTADO B.7.

- (a) Si \mathbf{A} y \mathbf{B} son ambas matrices de orden $m \times n$, entonces

$$\text{tr} \mathbf{A}^\top \mathbf{B} = \text{vec}^\top \mathbf{A} \text{vec} \mathbf{B},$$

- (b) Si \mathbf{A}, \mathbf{B} y \mathbf{C} son de órdenes adecuados, entonces

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

donde $\text{vec}^\top \mathbf{A} = (\text{vec} \mathbf{A})^\top$.

Finalmente, tenemos el siguiente resultado

RESULTADO B.8. Sean $\mathbf{A}, \mathbf{B}, \mathbf{C}$ y \mathbf{D} matrices, tal que, el producto \mathbf{ABCD} está definido y es cuadrado, entonces

$$\text{tr } \mathbf{ABCD} = \text{vec}^\top \mathbf{D}^\top (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec } \mathbf{B} = \text{vec}^\top \mathbf{D} (\mathbf{A} \otimes \mathbf{C}^\top) \text{vec } \mathbf{B}^\top.$$

Sea $\mathbf{F} : S \rightarrow \mathbb{R}^{m \times p}$, $S \subset \mathbb{R}^{n \times q}$ una función matricial, podemos notar que

$$\text{vec } \mathbf{F}(\mathbf{X}) = \mathbf{f}(\text{vec } \mathbf{X})$$

esto permite obtener el diferencial de una función matricial considerando la relación

$$\text{vec } d\mathbf{F}(\mathbf{C}; \mathbf{U}) = d\mathbf{f}(\text{vec } \mathbf{C}; \text{vec } \mathbf{U})$$

en cuyo caso \mathbf{F} tiene matriz Jacobiana

$$D\mathbf{F}(\mathbf{C}) = D\mathbf{f}(\text{vec } \mathbf{C})$$

Las consideraciones anteriores motivan el primer teorema de indentificación para funciones matriciales (Magnus y Neudecker, 1985)

TEOREMA B.9 (Primer teorema de indentificación para funciones matriciales). Sea $\mathbf{F} : S \rightarrow \mathbb{R}^{m \times p}$, $S \subset \mathbb{R}^{n \times q}$ función diferenciable, $\mathbf{C} \in S$ y \mathbf{U} matriz $n \times q$. Entonces

$$\text{vec } d\mathbf{F}(\mathbf{C}; \mathbf{U}) = (D\mathbf{F}(\mathbf{C})) \text{vec } \mathbf{U}.$$

con $(D\mathbf{F}(\mathbf{C}))^\top$ la matriz gradiente de \mathbf{F} .

B.3. Matriz Hessiana

Considere $\phi : S \rightarrow \mathbb{R}$ con $S \subset \mathbb{R}^n$, entonces se define la *matriz Hessiana* como la matriz de segundas derivadas, dada por

$$\mathbf{H}\phi(\mathbf{x}) = \frac{\partial^2 \phi(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^\top} = \frac{\partial}{\partial \mathbf{x}^\top} \left(\frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} \right)^\top = D(D\phi(\mathbf{x}))^\top.$$

Es posible definir el diferencial de funciones vectoriales y matriciales de manera análoga a la delineada anteriormente. Sin embargo, en este apéndice nos enfocaremos solamente en el cálculo de diferenciales de funciones escalares. El segundo diferencial de una función escalar está dado por

$$d^2 \phi = d(d\phi).$$

Magnus y Neudecker (1985) enunciaron el siguiente teorema de indentificación para matrices Hessianas de funciones escalares

TEOREMA B.10 (Segundo teorema de indentificación). Sea $\phi : S \rightarrow \mathbb{R}$, $S \subset \mathbb{R}^n$ dos veces diferenciable, $\mathbf{c} \in S$ y \mathbf{u} vector n -dimensional. Entonces

$$d^2 \phi(\mathbf{c}; \mathbf{u}) = \mathbf{u}^\top (\mathbf{H}\phi(\mathbf{c})) \mathbf{u}.$$

donde $\mathbf{H}\phi(\mathbf{c}) \in \mathbb{R}^{n \times n}$ es la matriz Hessiana de ϕ .

Algunas ventajas (prácticas) importantes del cálculo de diferenciales son:

- Sea $\mathbf{f}(\mathbf{x})$ función vectorial $m \times 1$ con argumento \mathbf{x} , vector n -dimensional, entonces

$$D\mathbf{f}(\mathbf{x}) \in \mathbb{R}^{m \times n} \quad \text{sin embargo,} \quad d\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$$

- Para funciones matriciales, $d\mathbf{F}(\mathbf{X})$ tiene la misma dimensión que \mathbf{F} sin importar la dimensión de \mathbf{X} .

B.4. Reglas fundamentales

A continuación se presentan algunas reglas fundamentales para el cálculo de diferenciales

Considere u y v funciones escalares y α una constante, entonces:

$$\begin{aligned} d\alpha &= 0, & d(\alpha u) &= \alpha du, & d(u+v) &= du + dv, \\ d(uv) &= (du)v + u(dv) & d(u/v) &= \frac{(du)v - u(dv)}{v^2}, (v \neq 0), \\ du^\alpha &= \alpha u^{\alpha-1} du, & de^u &= e^u du, \\ d\log u &= u^{-1} du, (u > 0) & d\alpha^u &= \alpha^u \log \alpha du, (\alpha > 0), \end{aligned}$$

aquí por ejemplo,

$$\phi(x) = u(x) + v(x).$$

Análogamente para U, V funciones matriciales, α un escalar (constante) y $A \in \mathbb{R}^{m \times n}$ constante, tenemos

$$\begin{aligned} dA &= 0, & d(\alpha U) &= \alpha dU, \\ d(U+V) &= dU + dV, & d(UV) &= (dU)V + U dV, \\ d(U \otimes V) &= dU \otimes dV, & d(U \odot V) &= dU \odot dV, \\ dU^\top &= (dU)^\top, & d\text{vec } U &= \text{vec } dU, & d\text{tr } U &= \text{tr } dU. \end{aligned}$$

Otros diferenciales de uso frecuente en Estadística son:

$$\begin{aligned} d|F| &= |F| \text{tr } F^{-1} dF, & d\log |F| &= \text{tr } F^{-1} dF, \\ dF^{-1} &= -F^{-1}(dF)F^{-1}. \end{aligned}$$

Bibliografia

- Andrews, D.F., Mallows, C.L. (1974). Scale mixtures of normal distributions. *Journal of the Royal Statistical Society, Series B* **36**, 99-102.
- Arellano, R. (1994). *Distribuições Elípticas: Propriedades, Inferência e Aplicações a Modelos de Regressão*. (Unpublished doctoral dissertation). Department of Statistics, University of São Paulo, Brazil.
- Atkinson, A.C. (1981). Two graphical displays for outlying and influential observations in regression. *Biometrika* **68**, 13-20.
- Atkinson, A.C. (1985). *Plots, Transformations and Regressions*. Oxford University Press, Oxford.
- Barlow, J.L. (1993). Numerical aspects of solving linear least squares problems. En *Handbook of Statistics, Vol. 9*, C.R. Rao (Ed.). Elsevier, Amsterdam, pp. 303-376.
- Barrodale, I., Roberts, F.D.K. (1973). An improved algorithm for discrete L1 linear approximations. *SIAM Journal of Numerical Analysis* **10**, 839-848.
- Barrodale, I., Roberts, F.D.K. (1974). Solution of an overdetermined system of equations in the L1 norm. *Communications of the ACM* **17**, 319-320.
- Belsley, D.A., Kuh, E., Welsh, R.E. (1980). *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. Wiley, New York.
- Björck, A. (1996). *Numerical Methods for Least Squares Problems*. Society for Industrial and Applied Mathematics, Philadelphia.
- Box, G.E.P., Cox, D.R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society, Series B* **26**, 211-252.
- Carroll, R.J., Ruppert, D. (1988). *Transformation and Weighting in Regression*. Chapman and Hall, New York.
- Charnes, A., Cooper, W.W., Ferguson, R.O. (1955). Optimal estimation of executive compensation by linear programming. *Management Science* **1**, 138-151.
- Chatterjee, S., Hadi, A.S. (1988). *Sensitivity Analysis in Linear Regression*. Wiley, New York.
- Christensen, R. (2011). *Plane Answers to Complex Questions: The Theory of Linear Models*, 4th Ed. Springer, New York.
- Cook, R.D. (1977). Detection of influential observation in linear regression. *Technometrics* **19**, 15-18.
- Cook, R.D., Weisberg, S. (1980). Characterizations of an empirical influence function for detecting influential cases in regression. *Technometrics* **22**, 495-508.
- Cook, R.D., Weisberg, S. (1982). *Residuals and Influence in Regression*. Chapman & Hall, New York.
- Daniel, C., Wood, F.S. (1980). *Fitting Equations to Data: Computer Analysis of Multifactor Data*, 2nd Ed. Wiley, New York.
- Davidian, M. Carroll, R.J. (1987). Variance function estimation. *Journal of the American Statistical Association* **82**, 1079-1091.

- Dempster, A.P., Laird, N.M., Rubin, D.B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B* **39**, 1-38.
- Díaz-García, J.A., Gutiérrez-Jáimez, R. (1999). *Cálculo Diferencial Matricial y Momentos de Matrices Aleatorias Elípticas*. Universidad de Granada.
- Dobson, A.J. (2002). *An Introduction to Generalized Linear Models*, 2nd Ed. Chapman & Hall, Boca Raton.
- Fahrmeir, L., Kneib, T., Lang, S., Marx, B. (2013). *Regression: Models, Methods and Applications*. Springer, Berlin.
- Fang, K.T., Kotz, S., Ng, K.W. (1990). *Symmetric Multivariate and Related Distributions*. Chapman & Hall, London.
- Galea, M. (1990). Técnicas de diagnóstico en regresión lineal. *Revista de la Sociedad Chilena de Estadística* **7**, 23-44.
- Gentle, J.E. (2007). *Matrix Algebra: Theory, Computation and Applications in Statistics*. Springer, New York.
- 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.
- Golub, G.H., Heath, M., Wahba, G. (1979). Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics* **21**, 215-223.
- Goodnight, J.H. (1979). A tutorial on the SWEEP operator. *The American Statistician* **33**, 149-158.
- Gorman, J.W., Toman, R.J. (1966). Selection of variables for fitting equations to data. *Technometrics* **8**, 27-51.
- Graybill, F.A. (1961). *An Introduction to Linear Statistical Models*. McGraw-Hill, New York.
- Graybill, F.A. (1976). *Theory and Application of the Linear Model*. Wadsworth & Brooks, Pacific Grove, CA.
- Graybill, F.A. (1983). *Matrices with Applications in Statistics*, 2nd Ed. Wadsworth, Belmont, CA.
- Groß, J. (2003). *Linear Regression*. Springer, Berlin.
- Gruber, M.H.J. (1998). *Improving Efficiency by Shrinkage*. Marcel Dekker, New York.
- Hald, A. (1952). *Statistical Theory with Engineering Applications*. Wiley, New York.
- Harville, D.A. (1997). *Matrix Algebra from a Statistician's Perspective*. Springer, New York.
- Hoaglin, D.C., Welsch, R.E. (1978). The hat matrix in regression and ANOVA. *The American Statistician* **32**, 17-22.
- Hocking, R. (1996). *Methods and Applications of Linear Models*. Wiley, New York.
- Hoerl, A.E., Kennard, R.W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* **12**, 55-67.
- Hoerl, A.E., Kennard, R.W. (1970). Ridge regression: Applications to nonorthogonal problems. *Technometrics* **12**, 69-82.
- Hoerl, A.E., Kennard, R.W., Baldwin, K.F. (1975). Ridge regression: Some simulations. *Communications in Statistics: Theory and Methods* **4**, 105-123.
- Huber, P.J. (1981). *Robust Statistics*. Wiley, New York.
- Kariya, T., Kurata, H. (2004). *Generalized Least Squares*. Wiley, Chichester.
- Lange, K. (1999). *Numerical Analysis for Statisticians*. Springer, New York.

- Lange, K.L., 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.
- Lange, K., Sinsheimer, J.S. (1993). Normal/independent distributions and their applications in robust regression. *Journal of Computational and Graphical Statistics* **2**, 175-198.
- Lawless, J.F., Wang, P. (1976). A simulation study of ridge and other regression estimators. *Communications in Statistics: Theory and Methods* **5**, 307-323.
- Lindley, D.V., Smith, A.F.M. (1972). Bayes estimates for the linear model *Journal of the Royal Statistical Society, Series B* **34**, 1-41.
- Little, R.J.A. (1988). Robust estimation of the mean and covariance matrix from data with missing values. *Applied Statistics* **37**, 23-38.
- Louis, T.A. (1982). Finding the observed information matrix when using the EM algorithm. *Journal of the Royal Statistical Society, Series B* **44**, 226-233.
- Magnus, J.R., Neudecker, H. (1985). Matrix differential calculus with applications to simple, Hadamard and Kronecker products. *Journal of Mathematical Psychology* **29**, 474-492.
- Magnus, J.R., Neudecker, H. (2007). *Matrix Differential Calculus with Applications in Statistics and Econometrics*, 3rd Ed. Wiley, New York.
- Magnus, J.R. (2010). On the concept of matrix derivative. *Journal of Multivariate Analysis* **101**, 2200-2206.
- McIntosh, A. (1982). *Fitting Linear Models: An Application of Conjugate Gradients Algorithms*. Springer, New York.
- McLachlan, G.J., Krishnan, T. (2008). *The EM Algorithm and Extensions*, 2nd Ed. Wiley, New York.
- O'Leary, D.P. (1990). Robust regression computation using iteratively reweighted least squares. *SIAM Journal on Matrix Analysis and Applications* **11**, 466-480.
- Osborne, M.R. (1985). *Finite Algorithms in Optimization and Data Analysis*. Wiley, New York.
- Osorio, F., Ogeda, A. (2021). *fastmatrix: Fast computation of some matrices useful in statistics*. R package version 0.3-819. URL: faosorios.github.io/fastmatrix
- Paula, G.A. (2013). *Modelos de Regressão, com Apoio Computacional*. Instituto de Matemática e Estatística, Universidade de São Paulo, Brasil.
- R Core Team (2020). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL: www.R-project.org
- Rao, C.R., Toutenburg, H., Shalabh, Heumann, C. (2008). *Linear Models and Generalizations: Least Squares and Alternatives*. Springer, New York.
- Ravishanker, N., Dey, D.K. (2002). *A First Course in Linear Model Theory*. Chapman & Hall, London.
- Ruppert, D., Wand, M.P., Carroll, R.J. (2003). *Semiparametric Regression*. Cambridge University Press, Cambridge.
- Schlossmacher, E.J. (1973). An iterative technique for absolute deviations curve fitting. *Journal of the American Statistical Association* **68**, 857-859.
- Searle, S.R. (1971). *Linear Models*. Wiley, New York.
- Searle, S.R. (1982). *Matrix Algebra Useful for Statistics*. Wiley, New York.
- Seber, G.A.F., Lee, A.J. (2003). *Linear Regression Analysis*, 2nd Ed. Wiley, New York.
- Sen, A., Srivastava, M. (1990). *Regression Analysis: Theory, Methods and Applications*. Springer, New York.

- Tong, Y.L. (1990). *The Multivariate Normal Distribution*. Springer, New York.
- Velleman, P.F., Welsch, R.E. (1981). Efficient computing of regression diagnostics. *The American Statistician* **35**, 234-242.
- Venables, W.N., Ripley, B.D. (2002). *Modern Applied Statistics with S*, 4th Ed. Springer, New York.
- Weisberg, S. (2005). *Applied Linear Regression*, 3rd Ed. Wiley, New York
- Welsch, R.E., Kuh, E. (1977). Linear regression diagnostics. *Massachusetts Institute of Technology and National Bureau of Economics, Cambridge*. Working paper No. 173.
- Woods, H., Steinour, H.H., Starke, H.R. (1932). Effect of composition of Portland cement on heat evolved during hardening. *Industrial and Engineering Chemistry* **24**, 1207-1214.
- Wu, C.F.J. (1983). On the convergence properties of the EM algorithm. *The Annals of Statistics* **11**, 95-103.