

**TEMA 8: Introducción a la teoría general de modelos lineales: regresión y análisis de la varianza**

- Descripción del modelo lineal general. Modelo de Gauss-Markov.
- Estimación de un modelo de Gauss-Markov.
- Inferencia bajo hipótesis de normalidad.
- Modelo de regresión lineal simple.
- Análisis de la varianza de una vía.

## I. El modelo lineal general. Modelo de Gauss-Markov

### MODELO LINEAL GENERAL

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (\mathbf{Y} = \mathbf{x}_1\beta_1 + \cdots, \mathbf{x}_k\beta_k + \boldsymbol{\varepsilon})$$

- $\mathbf{Y} = (Y_1, \dots, Y_n)^T$  es un vector aleatorio  $n$ -dimensional *observable*.
- $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$  es una *matriz conocida* de dimensión  $n \times k$  ( $k < n$ ), denominada *matriz de diseño*:

$$\mathbf{x}_j = (x_{1j}, \dots, x_{nj})^T, \quad j = 1, \dots, k$$

El rango de  $\mathbf{X}$  determina el *rango del modelo*; si el rango es  $k$ , el modelo es de *rango máximo o completo*.

- $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)^T$  es un vector de *parámetros desconocidos*, denominado *vector de efectos*, cuyas componentes ponderan los efectos de los vectores columna de  $\mathbf{X}$  en el vector  $\mathbf{Y}$ .
- $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$  es un vector aleatorio *no observable*, llamado *vector de errores*, que representa el error que se comete si se describe  $\mathbf{Y}$  por  $\mathbf{X}\boldsymbol{\beta}$ .



### MODELO DE GAUSS-MARKOV

*Es un modelo lineal en el que las componentes del vector de errores son variables aleatorias de segundo orden, centradas, homocedásticas (igual varianza) e incorreladas:*

$$E[\varepsilon_i] = 0, \quad \text{Var}[\varepsilon_i] = E[\varepsilon_i^2] = \sigma^2, \quad i = 1, \dots, n, \quad \text{Cov}[\varepsilon_i, \varepsilon_j] = E[\varepsilon_i \varepsilon_j] = 0, \quad i \neq j = 1, \dots, n.$$

*o, equivalentemente:*

$$E[\boldsymbol{\varepsilon}] = 0, \quad \text{Cov}[\boldsymbol{\varepsilon}] = E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T] = \sigma^2 \mathbf{I}_{n \times n},$$



*Objetivo*

Realizar inferencia sobre los parámetros  $\beta_1, \dots, \beta_k, \sigma^2$  a partir de una observación del vector  $\mathbf{Y}$ .

### Modelo de Gauss-Markov

$$\begin{aligned}
 \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} &\longrightarrow \begin{cases} E[\boldsymbol{\varepsilon}] = 0, & \text{Cov}[\boldsymbol{\varepsilon}] = E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T] = \sigma^2 \mathbf{I}_{n \times n} \\ E[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}, & \text{Cov}[\mathbf{Y}] = \sigma^2 \mathbf{I}_{n \times n}. \end{cases} \\
 Y_i = \sum_{j=1}^k x_{ij}\beta_j + \varepsilon_i, \quad i = 1, \dots, n &\longrightarrow \begin{cases} E[\varepsilon_i] = 0, & \text{Var}[\varepsilon_i] = \sigma^2, & \text{Cov}[\varepsilon_i, \varepsilon_j] = 0, \quad i \neq j \\ E[Y_i] = \sum_{j=1}^k x_{ij}\beta_j, & \text{Var}[Y_i] = \sigma^2, & \text{Cov}[Y_i, Y_j] = 0, \quad i \neq j. \end{cases}
 \end{aligned}$$

### I a) ESTIMACIÓN DE MÍNIMOS CUADRADOS DEL VECTOR DE EFECTOS

$$\begin{aligned}
 \text{Minimizar } S^2(\boldsymbol{\beta}) &= \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n \left( Y_i - \sum_{j=1}^k x_{ij}\beta_j \right)^2 = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\
 \bullet \frac{\partial S^2(\boldsymbol{\beta})}{\partial \beta_h} &= -2 \sum_{i=1}^n \left( Y_i - \sum_{j=1}^k x_{ij}\beta_j \right) x_{ih} = 0, \quad h = 1, \dots, k.
 \end{aligned}$$

**Ecuaciones normales:**  $\sum_{i=1}^n Y_i x_{ih} = \sum_{i=1}^n \sum_{j=1}^k x_{ij} x_{ih} \beta_j, \quad h = 1, \dots, k \longrightarrow \mathbf{X}^T \mathbf{Y} = (\mathbf{X}^T \mathbf{X}) \boldsymbol{\beta}.$

↓

**Estimador de mínimos cuadrados de  $\boldsymbol{\beta}$**   $\rightarrow \hat{\boldsymbol{\beta}}(\mathbf{Y}) = \left( \hat{\beta}_1(\mathbf{Y}), \dots, \hat{\beta}_k(\mathbf{Y}) \right)^T.$

- **Existencia:** existe, al menos, un estimador de mínimos cuadrados de  $\boldsymbol{\beta}$ .
- **Unicidad:** no está garantizada, salvo si *el modelo es de rango máximo*:

$$\hat{\boldsymbol{\beta}} := \hat{\boldsymbol{\beta}}(\mathbf{Y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \rightarrow \text{Función lineal de } \mathbf{Y} \rightarrow \begin{cases} E[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}. \\ \text{Cov}[\hat{\boldsymbol{\beta}}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \end{cases}$$

**Función estimable (escalar):**  $\psi(\boldsymbol{\beta})$  es estimable si admite un estimador insesgado, función lineal de las componentes de  $\mathbf{Y}$ :

$$\psi(\boldsymbol{\beta}) \text{ estimable} \Leftrightarrow \exists \mathbf{c} \in \mathbb{R}^n / E[\mathbf{c}^T \mathbf{Y}] = \psi(\boldsymbol{\beta}), \quad \forall \boldsymbol{\beta},$$

o, equivalentemente,  $\psi(\boldsymbol{\beta}) = \mathbf{c}^T \mathbf{X}\boldsymbol{\beta}, \quad \mathbf{c} \in \mathbb{R}^n.$

- Si el modelo es de rango máximo,  $\psi(\boldsymbol{\beta})$  es estimable  $\Leftrightarrow \psi(\boldsymbol{\beta}) = \mathbf{a}^T \boldsymbol{\beta}, \quad \mathbf{a} \in \mathbb{R}^k.$

**Teorema de Gauss-Markov:** Si  $\mathbf{a}^T \boldsymbol{\beta}$  es estimable, admite un único estimador lineal insesgado uniformemente de mínima varianza en la clase de estimadores lineales insesgados. Dicho estimador es  $\mathbf{a}^T \hat{\boldsymbol{\beta}}(\mathbf{Y})$ , y se denomina estimador de mínimos cuadrados de  $\mathbf{a}^T \boldsymbol{\beta}$ .

## I b) MODELO ESTIMADO, RESIDUOS Y ESTIMACIÓN DE LA VARIANZA

$$\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_k)^T := \hat{\beta}(\mathbf{Y}) \text{ estimador de mínimos cuadrados de } \beta$$

$$\text{Modelo estimado: } \hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} \longrightarrow \hat{Y}_i = \sum_{j=1}^k x_{ij}\hat{\beta}_j, \quad i = 1, \dots, n.$$

$$\text{Residuos mínimo-cuadráticos: } \mathbf{R} = \mathbf{Y} - \mathbf{X}\hat{\beta} \longrightarrow R_i = Y_i - \hat{Y}_i, \quad i = 1, \dots, n.$$

### Propiedades:

- $\hat{Y}_i$  es el estimador lineal insesgado de mínima varianza (estimador de mínimos cuadrados) de  $E[Y_i] = \sum_{j=1}^k x_{ij}\beta_j$ ,  $\forall i = 1, \dots, n$ .
- Los residuos son variables aleatorias con media nula,  $E[R_i] = 0$ ,  $\forall i = 1, \dots, n$ .
- El vector de residuos es ortogonal a los vectores columna de  $\mathbf{X}$ ,  $\mathbf{X}^T \mathbf{R} = 0$ .<sup>1</sup>
- El vector de residuos es ortogonal al vector estimado,  $\hat{\mathbf{Y}}^T \mathbf{R} = 0$ .

$$\text{VARIANZA RESIDUAL: } S_R^2 = \frac{\sum_{i=1}^n R_i^2}{n-r} = \frac{\|\mathbf{R}\|^2}{n-r} = \frac{\|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2}{n-r} \longrightarrow \text{estimador insesgado de } \sigma^2$$

<sup>1</sup>Esto indica que existen  $k$  relaciones lineales entre los residuos  $R_1, \dots, R_n$ , determinadas por las  $k$  columnas de  $\mathbf{X}$ ,  $\mathbf{x}_j^T \mathbf{R} = \sum_{i=1}^n x_{ij}R_i = 0$ ,  $j = 1, \dots, k$ . Por tanto, si  $\mathbf{X}$  es de rango  $r$ , el número de residuos  $R_i$  linealmente independientes es  $n - r \longrightarrow$  Los residuos tienen  $n - r$  grados de libertad.

**I c) INFERENCIA BAJO HIPÓTESIS DE NORMALIDAD (Ver Apéndice)**

$$\begin{aligned} \mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ \boldsymbol{\varepsilon} &\rightarrow \mathcal{N}_n(0, \sigma^2 I_{n \times n}) \end{aligned} \Leftrightarrow \mathbf{Y} \rightarrow \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 I_{n \times n})$$

► **Estimadores de máxima verosimilitud:**

*Función de verosimilitud:*  $\mathbf{y} \in \mathbb{R}^n \rightarrow L_{\mathbf{y}}(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left\{ -\frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2}{2\sigma^2} \right\}.$

■ *Estimador máximo verosímil de  $\boldsymbol{\beta} \rightarrow \hat{\boldsymbol{\beta}}$  (mínimos cuadrados).<sup>2</sup>*

■ *Estimador máximo verosímil de  $\sigma^2 \rightarrow \hat{\sigma}^2 = \frac{\sum_{i=1}^n R_i^2}{n} = \frac{(n-r)S_R^2}{n}$ .<sup>3</sup>*

► **Test de razón de verosimilitudes para la hipótesis lineal general:**

*Hipótesis lineal general:*  $H_0 : \mathbf{C}\boldsymbol{\beta} = \mathbf{0}$ , siendo  $\mathbf{C}_{q \times k}$  una matriz conocida de rango  $q$  ( $\leq k$ ), tal que todas las componentes del vector  $\mathbf{C}\boldsymbol{\beta}$  son estimables.

*Test de razón de verosimilitudes de tamaño  $\alpha$*

$$\varphi(\mathbf{Y}) = \begin{cases} 1 & F(\mathbf{Y}) > F_{q, n-r; \alpha} \\ 0 & F(\mathbf{Y}) \leq F_{q, n-r; \alpha} \end{cases} \quad F(\mathbf{Y}) = \frac{n-r}{q} \left( \frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}^0\|^2 - \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2} \right)$$

$\hat{\boldsymbol{\beta}}^0$ : estimador máximo verosímil de  $\boldsymbol{\beta}$  bajo  $H_0$

<sup>2</sup> Maximizar  $L_{\mathbf{y}}(\boldsymbol{\beta}, \sigma^2)$  en  $\boldsymbol{\beta}$  equivale a minimizar  $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ .

<sup>3</sup>  $\frac{\partial \ln L_{\mathbf{y}}(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = 0 \Rightarrow \hat{\sigma}^2(\mathbf{y}) = \frac{\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{n} = \frac{\sum_{i=1}^n r_i^2}{n}.$

## II. Modelo de regresión lineal simple

*Hipótesis:*  $X, Y$  variables aleatorias tales que  $E[Y^2] < +\infty$  y, fijado un valor arbitrario,  $X = x$ :

$$E[Y/X = x] = \beta_0 + \beta_1 x, \quad \text{Var}[Y/X = x] = \sigma^2.$$

$\Downarrow$

*Formulación:*  $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1 \dots n.$

- $x_i, \quad i = 1 \dots n$ , valores arbitrarios (fijos) de  $X$  (al menos dos distintos).
- $Y_i, \quad i = 1 \dots n$ , variables aleatorias que describen las observaciones de  $Y$ , supuesto que  $X = x_i$  ( $Y_i \equiv Y/X = x_i$ ).
- $\varepsilon_i, \quad i = 1 \dots n$ , variables aleatorias tales que  $E[\varepsilon_i] = 0, \text{Var}[\varepsilon_i] = \sigma^2, \quad i = 1, \dots, n.$

$\Downarrow \quad Y_1, \dots, Y_n$  independientes

**Modelo de Gauss-Markov de rango máximo (2).**

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}_{n \times 1} \quad \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}_{n \times 2} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}_{2 \times 1} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}_{n \times 1}$$

## II a) ESTIMACIÓN DEL MODELO <sup>4</sup>

### ► Estimador de mínimos cuadrados de $\beta$ :

Modelo de rango máximo  $\rightarrow \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$  <sup>5</sup>  $\left\{ \begin{array}{l} \text{Único y función lineal de } \mathbf{Y}. \\ E[\hat{\beta}] = \beta, \quad Cov[\hat{\beta}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \end{array} \right.$

$$\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)^T \rightarrow \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = \frac{\sigma_{xY}}{\sigma_x^2}$$

$\hat{\beta}_0$  y  $\hat{\beta}_1$  son los estimadores de mínimos cuadrados de  $\beta_0$  y  $\beta_1$ , respectivamente.

↓

Son los de mínima varianza uniformemente en la clase de estimadores lineales insesgados:

$$\star E[\hat{\beta}_0] = \beta_0, \quad E[\hat{\beta}_1] = \beta_1.$$

$$\star Var[\hat{\beta}_0] = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{n\sigma_x^2} \right), \quad Var[\hat{\beta}_1] = \sigma^2 \frac{1}{n\sigma_x^2}, \quad Cov[\hat{\beta}_0, \hat{\beta}_1] = -\sigma^2 \frac{\bar{x}}{n\sigma_x^2}.$$

### ► Modelo estimado, residuos mínimo cuadráticos y propiedades:

- $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} \rightarrow \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \rightarrow \hat{Y}_i = \bar{Y} + \frac{\sigma_{xY}}{\sigma_x^2} (x_i - \bar{x}), \quad i = 1, \dots, n.$
- $\mathbf{R} = \mathbf{Y} - \mathbf{X}\hat{\beta} \rightarrow R_i = Y_i - \hat{Y}_i \rightarrow R_i = Y_i - \bar{Y} - \frac{\sigma_{xY}}{\sigma_x^2} (x_i - \bar{x}), \quad i = 1, \dots, n.$
- $\mathbf{X}^T \mathbf{R} = 0 \Rightarrow \sum_{i=1}^n R_i = 0 \Rightarrow \sum_{i=1}^n R_i/n = 0 \Rightarrow \sum_{i=1}^n \hat{Y}_i/n = \bar{Y}.$
- $\hat{\mathbf{Y}}^T \mathbf{R} = \sum_{i=1}^n \hat{Y}_i R_i = 0.$

### ► Varianza residual: $S_R^2 = \frac{\sum_{i=1}^n R_i^2}{n-2}.$

$$\begin{aligned} \text{4 } \bar{x} &= \frac{\sum_{i=1}^n x_i}{n}, \quad \sigma_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}, \quad \bar{Y} = \frac{\sum_{i=1}^n Y_i}{n}, \quad \sigma_Y^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n}, \quad \sigma_{xY} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{n}. \\ \text{5 } (\mathbf{X}^T \mathbf{X})^{-1} &= \begin{pmatrix} n & n\bar{x} \\ n\bar{x} & \sum_{i=1}^n x_i^2 \end{pmatrix}^{-1} = \frac{1}{n\sigma_x^2} \begin{pmatrix} \sum_{i=1}^n x_i^2/n & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \quad (\sigma_x^2 \neq 0 \text{ ya que } \exists x_i \neq x_j). \\ (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} &= \frac{1}{n\sigma_x^2} \begin{pmatrix} \sum_{i=1}^n x_i^2/n & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \begin{pmatrix} n\bar{Y} \\ \sum_{i=1}^n x_i Y_i \end{pmatrix} = \begin{pmatrix} \bar{Y} - \bar{x} \frac{\sigma_{xY}}{\sigma_x^2} \\ \frac{\sigma_{xY}}{\sigma_x^2} \end{pmatrix}. \end{aligned}$$

## II b) ANÁLISIS DE LA BONDAD DEL MODELO ESTIMADO

### Descomposición de la variabilidad de las observaciones <sup>6</sup>

$$VT = VE + VNE$$

- $VT = \sum_{i=1}^n (Y_i - \bar{Y})^2 = n\sigma_Y^2 \rightarrow$  Variabilidad total (de  $Y_1, \dots, Y_n$ ).
- $VE = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \frac{n\sigma_{xY}^2}{\sigma_x^2} \rightarrow$  Variabilidad explicada por el modelo (de  $\hat{Y}_1, \dots, \hat{Y}_n$ ).
- $VNE = \sum_{i=1}^n R_i^2 \rightarrow$  Variabilidad no explicada por el modelo (de  $R_1, \dots, R_n$ ).

### Coeficiente de determinación lineal

$$R^2 = \frac{VE}{VT} = \frac{\sigma_{xY}^2}{\sigma_x^2 \sigma_Y^2}.$$



Proporción de variabilidad total explicada por el modelo

<sup>6</sup>  $Y_i = \hat{Y}_i + R_i, \quad i = 1, \dots, n.$

$\Downarrow$

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i + R_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n R_i^2 + 2 \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y}) R_i}_{0} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n R_i^2.$$



## II c) PREDICCIÓN A PARTIR DEL MODELO ESTIMADO

$$X = x_p \longrightarrow Y_p = \beta_0 + \beta_1 x_p + \varepsilon_p; \quad E[\varepsilon_p] = 0, \quad Var[\varepsilon_p] = \sigma^2$$

$\Downarrow$

$$\text{Predicción de } Y \text{ cuando } X = x_p \longrightarrow \hat{Y}_p = \hat{\beta}_0 + \hat{\beta}_1 x_p = \bar{Y} + \frac{\sigma_{xY}}{\sigma_x^2} (x_p - \bar{x}).$$

- $E[\hat{Y}_p] = \beta_0 + \beta_1 x_p = E[Y_p]$ .
- $Var[\hat{Y}_p] = E[(\hat{Y}_p - E[\hat{Y}_p])^2] = \sigma^2 \left[ \frac{1}{n} + \frac{(x_p - \bar{x})^2}{n\sigma_x^2} \right]$ .
- $ECM[\hat{Y}_p] = E[(\hat{Y}_p - Y_p)^2] = Var[\hat{Y}_p] + Var[Y_p] = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{n\sigma_x^2} \right]$ .

## II d) APLICACIÓN PRÁCTICA: RECTA DE REGRESIÓN ESTIMADA

$(x_1, y_1), \dots, (x_n, y_n)$  observaciones independientes del vector aleatorio  $(X, Y)$ .

*Modelo lineal correspondiente a  $x_1, \dots, x_n \rightarrow Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$ .*

*Modelo estimado a partir de  $Y_1, \dots, Y_n \rightarrow \hat{Y}_i = \bar{Y} + \frac{\sigma_{xY}}{\sigma_x^2}(x_i - \bar{x}), \quad i = 1, \dots, n$ .*

$$\downarrow Y_1 = y_1, \dots, Y_n = y_n$$

$$\hat{y}_i = \bar{y} + \frac{\sigma_{xy}}{\sigma_x^2}(x_i - \bar{x}), \quad i = 1, \dots, n.$$

*Recta de regresión estimada <sup>7</sup>a partir de  $(x_1, y_1), \dots, (x_n, y_n) \rightarrow$*

$$y = \bar{y} + \frac{\sigma_{xy}}{\sigma_x^2}(x - \bar{x})$$

*Coefficiente de determinación lineal estimado <sup>8</sup> a partir de  $(x_1, y_1), \dots, (x_n, y_n) \rightarrow$*

$$r^2 = \frac{\sigma_{xy}^2}{\sigma_x^2 \sigma_y^2}$$

*Predicción a partir de la recta de regresión estimada  $\rightarrow \hat{y}_p = \bar{y} + \frac{\sigma_{xy}}{\sigma_x^2}(x_p - \bar{x})$ .*

*Estimación del error cuadrático medio de la predicción  $\rightarrow \widehat{ECM}(\hat{Y}_p) = s_R^2 \left[ 1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{n\sigma_x^2} \right]$ .*

<sup>7</sup> Recta de regresión teórica:  $y = E[Y] + \frac{Cov[X, Y]}{Var[X]}(x - E[X])$

<sup>8</sup> Coeficiente de determinación lineal de  $X$  e  $Y$ :  $\rho_{XY}^2 = \frac{Cov^2[X, Y]}{Var[X]Var[Y]}$

## II e) CONTRASTE DE REGRESIÓN BAJO HIPÓTESIS DE NORMALIDAD

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1 \dots n, \quad \varepsilon_i \rightarrow \mathcal{N}(0, \sigma^2) \quad \left( \Leftrightarrow Y_i \rightarrow \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2) \right)$$

$Y_1, \dots, Y_n$  independientes.

$$H_0 : \beta_1 = 0 \quad \longleftrightarrow \quad H_0 : \mathbf{C}\boldsymbol{\beta} = 0 \text{ con } \mathbf{C} = (0, 1)_{1 \times 2} \quad (\mathbf{C}\boldsymbol{\beta} \text{ estimable y rango } q = 1).$$

*Estadístico de contraste del test de razón de verosimilitudes*

$$F(\mathbf{Y}) = \frac{n-r}{q} \left( \frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}^0\|^2 - \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2} \right) \rightarrow \begin{cases} \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 = \sum_{i=1}^n R_i^2 = VNE \\ \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}^0\|^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 = VT^{(*)} \end{cases}$$

(\*) Bajo  $H_0$ ,  $Y_i = \beta_0 + \varepsilon_i \rightarrow \mathcal{N}(\beta_0, \sigma^2)$ ,  $i = 1, \dots, n \Rightarrow (Y_1, \dots, Y_n)$  m.a.s. de  $\mathcal{N}(\beta_0, \sigma^2) \Rightarrow \hat{\boldsymbol{\beta}}_0^0 = \bar{Y}$ .

$$\Downarrow$$

$$\hat{\boldsymbol{\beta}}^0 = (\hat{\beta}_0^0, \hat{\beta}_1^0)^T = (\bar{Y}, 0)^T \Rightarrow \mathbf{X}\hat{\boldsymbol{\beta}}^0 = (\bar{Y}, \dots, \bar{Y})^T.$$



$$F(\mathbf{Y}) = (n-2) \frac{VT - VNE}{VNE} = \frac{VE}{VNE/(n-2)} = \frac{VE}{S_R^2} \xrightarrow{H_0} F(1, n-2).$$

Test de razón de verosimilitudes de tamaño  $\alpha \rightarrow \varphi(\mathbf{Y}) = \begin{cases} 1 & F(\mathbf{Y}) > F_{1, n-2; \alpha} \\ 0 & F(\mathbf{Y}) \leq F_{1, n-2; \alpha} \end{cases}$

**Nota:**  $F(\mathbf{Y}) = (n-2) \frac{R^2}{1-R^2} \rightarrow$  función creciente de  $R^2$ . Así, grandes valores de  $R^2$  conducen al rechazo de  $H_0$ , lo que concuerda con la definición de  $R^2$  como medida de bondad del modelo estimado.

### III. Análisis de la varianza de una vía

Técnica estadística (Fisher) para *determinar si un supuesto factor de variación afecta al comportamiento de una cierta variable aleatoria*. Se aplica bajo los siguientes supuestos:

- *La variable de interés no está afectada por factores distintos al que es objeto de estudio.*
- *El factor de variación tiene un número finito de niveles ( $k$ ) y, en cada uno de ellos, la variable tiene distribución normal, con la misma varianza:*

$Y_i$  : variable de interés en el nivel  $i$ -ésimo  $\rightarrow \mathcal{N}(\mu_i, \sigma^2)$ ,  $i = 1, \dots, k$ .

↓ Hipótesis a contrastar

*El supuesto factor de variación no afecta al comportamiento de la variable:*

$$H_0 : \mu_1 = \dots = \mu_k$$

*Igualdad de medias de  $k$  poblaciones normales con varianza común.*

**Resolución:**

- $(Y_{i1}, \dots, Y_{in_i})$  muestra aleatoria simple de  $Y_i$ ,  $i = 1, \dots, k$ , todas independientes.
- $Y_{ij} \rightarrow \mathcal{N}(\mu_i, \sigma^2) \Rightarrow \underbrace{Y_{ij} = \mu_i + \varepsilon_{ij}}_{i=1, \dots, k, j=1, \dots, n_i}$ , con  $\varepsilon_{ij} \rightarrow \mathcal{N}(0, \sigma^2)$ .

$\mathbf{Y} = \mathbf{X}\boldsymbol{\mu} + \boldsymbol{\varepsilon}$  modelo lineal  $n$ -dimensional ( $n = \sum_{i=1}^k n_i$ )

$$\mathbf{Y} = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ Y_{k1} \\ \vdots \\ Y_{kn_k} \end{pmatrix}_{n \times 1} \quad \mathbf{X} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 1 & \cdots & 0 \\ \vdots & & & \vdots \\ 0 & 0 & \cdots & 1 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}_{n \times k} \rightarrow \text{rango} = k, \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_k \end{pmatrix}_{k \times 1} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \\ \varepsilon_{21} \\ \vdots \\ \varepsilon_{2n_2} \\ \vdots \\ \varepsilon_{k1} \\ \vdots \\ \varepsilon_{kn_k} \end{pmatrix}_{n \times 1}$$

$$E[\varepsilon_{ij}] = 0, \quad i = 1, \dots, k, \quad j = 1, \dots, n_i; \quad E[\varepsilon_{ij}\varepsilon_{i'j'}] = 0, \quad i \neq i' \text{ ó } j \neq j'$$

↓

**Modelo de Gauss-Markov de rango máximo ( $k$ )**

### III a) ESTIMACIÓN DEL MODELO

► Estimador máximo verosímil (mínimos cuadrados) de  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_k)^T$ <sup>9</sup>

$$\boxed{\hat{\mu}_i = \bar{Y}_i, \quad i = 1, \dots, k \longrightarrow \hat{\boldsymbol{\mu}} = (\bar{Y}_1, \dots, \bar{Y}_k)^T}$$

► Modelo estimado, residuos mínimo cuadráticos y propiedades:

- $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\mu}} = \left(\bar{Y}_1, \dots, \bar{Y}_1 \mid \dots \mid \bar{Y}_k, \dots, \bar{Y}_k\right)^T \longrightarrow \hat{Y}_{ij} = \bar{Y}_i, \quad i = 1, \dots, k, \quad j = 1, \dots, n_i.$
- $\mathbf{R} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\mu}} \longrightarrow R_{ij} = Y_{ij} - \hat{Y}_{ij} = Y_{ij} - \bar{Y}_i, \quad i = 1, \dots, k, \quad j = 1, \dots, n_i.$
- $\mathbf{X}^T \mathbf{R} = 0 \Rightarrow \sum_{j=1}^{n_i} R_{ij} = 0, \quad i = 1, \dots, k \Rightarrow \left\{ \begin{array}{l} \sum_{i=1}^k \sum_{j=1}^{n_i} R_{ij} / n = 0 \\ \sum_{i=1}^k \sum_{j=1}^{n_i} \hat{Y}_{ij} / n = \sum_{i=1}^k \sum_{j=1}^{n_i} \bar{Y}_i / n = \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij} / n = \bar{Y}. \end{array} \right.$
- $\hat{\mathbf{Y}}^T \mathbf{R} = \sum_{i=1}^k \sum_{j=1}^{n_i} \hat{Y}_{ij} R_{ij} = \sum_{i=1}^k \sum_{j=1}^{n_i} \bar{Y}_i R_{ij} = 0.$

---

<sup>9</sup>  $S^2(\boldsymbol{\mu}) = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \mu_i)^2 \Rightarrow \frac{\partial S^2(\boldsymbol{\mu})}{\partial \mu_h} = \frac{d \sum_{j=1}^{n_h} (Y_{hj} - \mu_h)^2}{d \mu_h} = -2 \sum_{j=1}^{n_h} (Y_{hj} - \mu_h) = 0 \Rightarrow \hat{\mu}_h = \frac{\sum_{j=1}^{n_h} Y_{hj}}{n_h} = \bar{Y}_h$

### III b) DESCOMPOSICIÓN DE LA VARIABILIDAD DE LAS OBSERVACIONES<sup>10</sup>

$$VT = VE + VNE$$

- $VT = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2 \rightarrow$  Variabilidad total (de  $Y_{ij}$ ,  $i = 1, \dots, k$ ,  $j = 1, \dots, n_i$ ).
- $VE = \sum_{i=1}^k \sum_{j=1}^{n_i} (\hat{Y}_{ij} - \bar{Y})^2 \rightarrow$  Variabilidad explicada por el modelo (de  $\hat{Y}_{ij}$ ,  $i = 1, \dots, k$ ,  $j = 1, \dots, n_i$ ).
- $VNE = \sum_{i=1}^k \sum_{j=1}^{n_i} R_{ij}^2 \rightarrow$  Variabilidad no explicada por el modelo (de  $R_{ij}$ ,  $i = 1, \dots, k$ ,  $j = 1, \dots, n_i$ ).

$$\downarrow \hat{Y}_{ij} = \bar{Y}_i, R_{ij} = Y_{ij} - \bar{Y}_i$$

- $VE = \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{Y}_i - \bar{Y})^2 \rightarrow$  Variabilidad de  $\bar{Y}_1, \dots, \bar{Y}_k \rightarrow$  Variabilidad entre grupos.
- $VNE = \sum_{i=1}^k \underbrace{\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2}_{\text{Variabilidad de las observaciones de la muestra } i\text{-ésima}} \rightarrow$  Variabilidad dentro de grupos.

$$^{10} Y_{ij} = \hat{Y}_{ij} + R_{ij}$$

$\downarrow$

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (\hat{Y}_{ij} + R_{ij} - \bar{Y})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (\hat{Y}_{ij} - \bar{Y})^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} R_{ij}^2 + 2 \underbrace{\sum_{i=1}^k \sum_{j=1}^{n_i} (\hat{Y}_{ij} - \bar{Y}) R_{ij}}_0$$

III c) PROBLEMA DE CONTRASTE  $H_0 : \mu_1 = \dots = \mu_k$

$$H_0 : \mu_1 - \mu_2 = 0, \mu_1 - \mu_3 = 0, \dots, \mu_1 - \mu_k = 0$$

$\updownarrow$

$$H_0 : \mathbf{C}\boldsymbol{\mu} = 0, \text{ con } \mathbf{C} = \begin{pmatrix} 1 & -1 & 0 & \dots & 0 \\ 1 & 0 & -1 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 0 & \dots & -1 \end{pmatrix}_{(k-1) \times k} \rightarrow \text{rango} = k - 1 \text{ (Hipótesis lineal general)}$$

*Estadístico de contraste del test de razón de verosimilitudes*

$$F(\mathbf{Y}) = \frac{n-r}{q} \left( \frac{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\mu}}^0\|^2}{\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\mu}}\|^2} \right) \rightarrow \begin{cases} \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\mu}}\|^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} R_{ij}^2 = VNE \\ \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\mu}}^0\|^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2 = VT^{(*)} \end{cases}$$

(\*) Bajo  $H_0$ ,  $Y_{ij} \rightarrow \mathcal{N}(\mu, \sigma^2)$ ,  $\forall i, j \Rightarrow (Y_{11}, \dots, Y_{1n_1}, \dots, Y_{k1}, \dots, Y_{kn_k})$  m.a.s. de  $\mathcal{N}(\mu, \sigma^2) \Rightarrow \hat{\mu} = \bar{Y}$ .

$$\downarrow \\ \hat{\mu}_i^0 = \bar{Y}, \quad i = 1, \dots, k \rightarrow \hat{\boldsymbol{\mu}}^0 = (\bar{Y}, \dots, \bar{Y})^T \rightarrow \mathbf{X}\hat{\boldsymbol{\mu}}^0 = (\bar{Y}, \dots, \bar{Y} | \dots | \bar{Y}, \dots, \bar{Y})^T.$$

$\Downarrow$

$$F(\mathbf{Y}) = \frac{n-k}{k-1} \frac{VT - VNE}{VNE} = \frac{VE/(k-1)}{VNE/(n-k)} \xrightarrow{H_0} F(k-1, n-k).$$

**TABLA ANOVA DE UNA VÍA**

Fuentes de Variación	Variabilidad	Grados de libertad	Varianzas
Entre grupos	$VE = \sum_{i=1}^k n_i (\bar{Y}_i - \bar{Y})^2$	$k - 1$	$S_E^2 = VE/(k - 1)$
Dentro de grupos	$VNE = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$	$n - k$	$S_R^2 = VNE/(n - k)$
Total	$VT = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y})^2$	$n - 1$	



**Test de razón de verosimilitudes de tamaño  $\alpha$**

$$\varphi(\mathbf{Y}) = \begin{cases} 1 & F(\mathbf{Y}) > F_{k-1, n-k; \alpha} \\ 0 & F(\mathbf{Y}) \leq F_{k-1, n-k; \alpha} \end{cases} \quad F(\mathbf{Y}) = \frac{S_E^2}{S_R^2}.$$

**Nota:** En la práctica, si el problema planteado no especifica un nivel de significación se trabaja con el denominado *p-nivel o p-valor asociado a los datos*:

$$p = P(F(k - 1, n - k) > F_{exp}),$$

siendo  $F_{exp}$  el valor del estadístico  $F(\mathbf{Y})$  obtenido de los datos concretos:

- Si el *p-nivel* es pequeño (usualmente, 0.05 o menor) se rechaza  $H_0$ .
- Si el *p-nivel* es grande (0.15 o mayor), se acepta  $H_0$ .
- Para valores intermedios hay que tratar cada situación en particular aunque, normalmente, es aconsejable tomar más datos y rehacer los cálculos.



## Apéndice: Distribución normal $n$ -dimensional

Un vector aleatorio  $Y = (Y_1, \dots, Y_n)^T$  tiene distribución normal ( $n$ -dimensional), lo que se denota  $Y \rightarrow \mathcal{N}_n(\mu, \Sigma)$ , si su función de densidad es de la forma:

$$f(y) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{(y - \mu)^T \Sigma^{-1} (y - \mu)}{2}}, \quad y \in \mathbb{R}^n,$$

donde  $\mu = (\mu_1, \dots, \mu_n)^T \in \mathbb{R}^n$  y  $\Sigma = ((\sigma_{ij}))_{n \times n}$  es una matriz definida positiva.

### Propiedades:

- $E[Y] = \mu$ ,  $Cov[Y] = E[(Y - \mu)(Y - \mu)^T] = \Sigma$ .
- Las distribuciones marginales de cualquier dimensión son normales y, en particular,  $Y_i \rightarrow \mathcal{N}(\mu_i, \sigma_{ii})$ ,  $i = 1, \dots, n$ .
- $Y \rightarrow \mathcal{N}_n(\mu, \Sigma) \Rightarrow \forall \gamma \in \mathbb{R}^n$  (vector constante),  $Y + \gamma \rightarrow \mathcal{N}_n(\mu + \gamma, \Sigma)$ .
- $Y = (Y_1, \dots, Y_n)^T \rightarrow \mathcal{N}_n(\mu, \Sigma)$ :  
 $Y_1, \dots, Y_n$  son independientes  $\Leftrightarrow \Sigma$  es diagonal  $\Leftrightarrow \rho_{Y_i, Y_j} = 0$ ,  $\forall i, j = 1, \dots, n$ ,  $i \neq j$ .