

## "MODELO LINEAL DE RANGO COMPLETO"

### Modelo matemático

$$y = g(x_1, x_2, \dots, x_k)$$

y                    variable a explicar  
x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>k</sub>    variables explicativas

De esta manera, conocidos los valores de las variables explicativas y la forma específica de g, es posible estimar el valor correspondiente de la variable y.

### Modelo estadístico

$$y = g(x_1, x_2, \dots, x_k) + \varepsilon$$

Incluye un término de error ε (aleatorio y no observable) que contempla todos los posibles factores generadores de discrepancias entre los valores observados y los valores estimados por el modelo: selección de variables, selección de individuos, selección de la ecuación, errores de diseño y errores de medición.

### Modelo lineal

$$y = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

Por lo general, esta ecuación lineal constituye una buena aproximación a la verdadera relación

Las constantes β<sub>1</sub>, β<sub>2</sub>, ... β<sub>k</sub>, se denominan parámetros, son desconocidas y se estiman a partir de información reportada por observaciones de las variables

Al efectuar n observaciones de las k+1 variables, las relaciones entre ellas quedan descritas mediante el conjunto de n ecuaciones:

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i \quad (i = 1, 2, \dots, n)$$

que en forma matricial pueden escribirse como:

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ \vdots \\ k \end{pmatrix} + \begin{pmatrix} 1 \\ 2 \\ \vdots \\ n \end{pmatrix}$$

o bien:

$$y_{nx1} = X_{nxk} \beta_{kx1} + \varepsilon_{nx1}$$

expresión que se conoce como modelo lineal general.

### Ajuste del modelo

Al hallarse un estimador  $\tilde{\beta}$  del vector de parámetros y sustituir en la ecuación anterior, se obtiene un ajuste:

$$\mathbf{y} = \mathbf{X}\tilde{\beta} + \mathbf{e}$$

surgiendo un término de error (aleatorio y observable) que se denomina residuo.

### Ajuste mínimo cuadrático

Es un procedimiento de estimación mediante el cual se obtiene el estimador  $\mathbf{b}$ , que minimiza la suma de cuadrados de los residuos:

$$\text{SCE} = \sum_{i=1}^n e_i^2 = \mathbf{e}^t \mathbf{e}$$

donde:

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\tilde{\beta}$$

Al derivar SCE respecto de  $\tilde{\beta}$  e igualando al vector nulo, se obtienen las ecuaciones normales:

$$\mathbf{X}^t \mathbf{X} \tilde{\beta} = \mathbf{X}^t \mathbf{y}$$

de donde se deduce el estimador mínimo cuadrático de  $\beta$ :

$$\mathbf{b} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}$$

asumiendo que  $\mathbf{X}_{n \times k}$  es una matriz de rango  $k$  ( $k < n$ ). (De ahí el nombre de modelo lineal de rango completo)

### Vector de estimaciones

$$\tilde{\mathbf{y}} = \mathbf{X}\mathbf{b}$$

resultando que:

$$\mathbf{e} = \mathbf{y} - \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{X}\mathbf{b} = (\mathbf{I} - \mathbf{H})\mathbf{y} = \mathbf{Q}\mathbf{y} = \mathbf{Q}\mathbf{e}$$

donde:

$\mathbf{H} = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t$  (matriz "hat") es una matriz simétrica, idempotente y de rango (o traza) igual a  $k$ , y

$\mathbf{Q} = \mathbf{I} - \mathbf{H}$  es una matriz simétrica, idempotente y de rango (o traza) igual a  $n-k$ .

Nota: A  $\mathbf{H}$  se denomina matriz hat ya que  $\mathbf{H}\mathbf{y} = \tilde{\mathbf{y}}$ .

### Varianza residual

$$s^2 = \frac{\text{SCE}}{n-k} = \frac{\mathbf{e}^t \mathbf{e}}{n-k} = \frac{t_{\mathbf{Q}}}{n-k}$$

### Modelo con término constante

En el caso en que se incluye un término constante  $\beta_0$ , el modelo estadístico queda

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_k + \varepsilon$$

siendo este modelo un caso particular del modelo lineal general en el cual:

$$\begin{aligned} \mathbf{y}_{n \times 1} &= \mathbf{j}_{n \times 1} \beta_0 + \mathbf{X}^* \boldsymbol{\beta}^* + \varepsilon \\ &= (\mathbf{j}_{n \times 1}, \mathbf{X}^*) \begin{pmatrix} \beta_0 \\ \boldsymbol{\beta}^* \end{pmatrix} + \varepsilon \\ &= \mathbf{X}_{n \times (k+1)} \boldsymbol{\beta}_{(k+1) \times 1} + \varepsilon \end{aligned}$$

es decir, en el que:

$$\mathbf{X}_{n \times (k+1)} = (\mathbf{j}_{n \times 1}, \mathbf{X}^*) \quad \text{y} \quad \boldsymbol{\beta}_{(k+1) \times 1} = \begin{pmatrix} \beta_0 \\ \boldsymbol{\beta}^* \end{pmatrix}$$

donde  $\mathbf{X}^*$  es la matriz que contiene la información relativa a las  $p$  variables explicativas, y  $\boldsymbol{\beta}$  el vector de coeficientes de regresión correspondientes:  $\beta_1, \beta_2, \dots, \beta_k$ . Resolviendo las ecuaciones normales para este caso, se obtiene:

$$\begin{aligned} b_0 &= \bar{y} - \bar{\mathbf{x}}_*^t \mathbf{b}^* \\ \mathbf{b}^* &= (\mathbf{X}_*^t \mathbf{X}_*)^{-1} \mathbf{X}_*^t \mathbf{y} = \mathbf{S}_{xx}^{-1} \mathbf{s}_{xy} \end{aligned}$$

El vector de estimaciones en este caso:

$$\tilde{\mathbf{y}} = \mathbf{j}_{n \times 1} b_0 + \mathbf{X}^* \mathbf{b}^*$$

resultando que:

$$\mathbf{e} = \mathbf{y} - \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{X}_*^* \mathbf{b}^* = (\mathbf{I} - \mathbf{H}_*) \mathbf{y} = \mathbf{Q}_* \mathbf{y} = \mathbf{Q}_* \mathbf{e}$$

donde:

$\mathbf{H}_* = \mathbf{X}_*^* (\mathbf{X}_*^* \mathbf{X}_*^*)^{-1} \mathbf{X}_*^{*t}$  es una matriz simétrica, idempotente y de rango (o traza) igual a  $k$ , y

$\mathbf{Q}_* = \mathbf{I} - \mathbf{H}_*$  es una matriz simétrica, idempotente y de rango (o traza) igual a  $n-k$ .

y la varianza residual:

$$s^2 = \frac{\text{SCE}}{n-k-1} = \frac{\mathbf{e}^t \mathbf{e}}{n-k-1} = \frac{\mathbf{y}^t \mathbf{Q}_* \mathbf{y}}{n-k-1}$$

**Modelo de regresión simple**(caso bidimensional, k=1)

$$\begin{aligned}b_0 &= \bar{y} - b_1 \bar{x} \\b_1 &= \frac{s(x, y)}{s^2(x)} = r(x, y) \frac{s(y)}{s(x)} \\s^2 &= \frac{SCE}{n - 2} = \frac{(n - 1)s^2(y)(1 - r^2)}{n - 2}\end{aligned}$$