

Capítulo 3: Regresión Múltiple

Fernando A. Crespo R.

14 de Octubre de 2021



Índice

3. Modelo Clásico de Regresión Múltiple Lineal

3.1 Matriz de Correlación

3.2 Marco Teórico

3.3 Supuestos

3.4 Estimadores del Modelo Clásico de Regresión Lineal

3.5 Coeficiente de Determinación

3.1 Matriz de Correlación

- ▶ Lo primero es estudiar la relación entre variables. Para ello se estudia la matriz de correlación. Donde se despliega los coeficientes de correlación para cada par de variables.
- ▶ Multicolinealidad: cuando dos variables presentan una alta correlación entre ellas. Ya que eso significa que hay una relación directa entre las dos variables. Por lo tanto, es difícil después, imputar cuanto aporta una variable para explicar la variación de otra. Dos variables con correlación alta no proporcionan información adicional.
- ▶ Regla para elegir:
 1. No debe haber correlación alta entre variables predictoras o explicativas (nosotros las denominamos como x).

3.1 Marco Teórico: Modelo Clásico de Regresión Múltiple Lineal

- ▶ El vector columna \mathbf{x}_j contiene las n observaciones con $k = 1, \dots, k$, añadamos el vector columna a la matriz $\mathbf{X} \in \mathbf{M}_{n \times k+1}(\mathbb{R})$, donde la primera columna es de unos ($\mathbf{x}_1 = \mathbf{1}_n$). Puede escribirse como:

$$\mathbf{y} = \mathbf{1}_n\beta_0 + \mathbf{x}_1\beta_1 + \mathbf{x}_2\beta_2 \dots + \mathbf{x}_k\beta_k + \epsilon, \quad (1)$$

o como:

$$\mathbf{y} = \beta_0 + \mathbf{x}_1\beta_1 + \mathbf{x}_2\beta_2 \dots + \mathbf{x}_k\beta_k + \epsilon. \quad (2)$$

ϵ es la perturbación aleatoria.

3.3 Supuestos

- ▶ Supuesto 1:

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon. \quad (3)$$

con

$$\mathbf{X} = [\mathbf{1}_n \ x_1 \ x_2 \ \dots \ x_k]. \quad (4)$$

- ▶ Supuesto 2: \mathbf{X} es una matriz $\in M_{n \times k+1}(IR)$ de rango completo, esta es la condición de identificación.
- ▶ Supuesto 3:

$$\mathbf{E}[\epsilon|\mathbf{X}] = \begin{bmatrix} E[\epsilon_1|X] \\ E[\epsilon_2|X] \\ \vdots \\ E[\epsilon_n|X] \end{bmatrix} = \mathbf{0}. \quad (5)$$

3.3 Supuestos

- ▶ Se pide además:

$$\mathbf{Var}[\epsilon_j | \mathbf{X}] = \sigma^2, \quad (6)$$

para $i = 1, \dots, n$, y

$$\mathbf{Cov}[\epsilon_i, \epsilon_j | \mathbf{X}] = \mathbf{0}, \quad (7)$$

para $i \neq j$.

- ▶ La varianza constante se denomina homocedasticidad, y en caso de que varíe por variable, se denomina heterocedasticidad.
(7) se denomina la no autocorrelación, o de forma equivalente, que el experimento se hizo con elección aleatoria de los datos.

3.3 Supuestos

- ▶ Supuesto 4:

$$\mathbb{E}[\epsilon\epsilon^t | \mathbf{X}] = \sigma^2 \mathbf{I}. \quad (8)$$

- ▶ Supuesto 5: \mathbf{X} es una matriz conocida de $M_{n \times k+1}(IR)$ de constantes (no es estocástica).
- ▶ Supuesto 6:

$$\epsilon | \mathbf{X} \sim \mathbf{N}[\mathbf{0}, \sigma^2 \mathbf{I}]. \quad (9)$$

3.4 Estimadores de la Regresión

- ▶ Los estimadores se pueden obtener por dos vías:
 - ▶ Obteniendo el Estimador Máximo Verosímil dado los supuestos sobre los errores, u
 - ▶ Obteniendo los mínimos cuadráticos.
- ▶ Vamos a obtener los estimadores a partir de minimizar:

$$\mathbf{S}(\beta) = (\mathbf{Y} - \mathbf{X}\beta)^t(\mathbf{Y} - \mathbf{X}\beta) \in \mathbb{R}. \quad (10)$$

- ▶ Desarrollando (10), obtenemos:

$$\mathbf{S}(\beta) = \mathbf{Y}^t\mathbf{Y} - 2\mathbf{Y}^t\mathbf{X}\beta + \beta^t(\mathbf{X}^t\mathbf{X})\beta. \quad (11)$$

- ▶ Para minimizar necesitamos $\mathbf{S}'(\beta)$ e igualamos a $\mathbf{0}$:

$$\mathbf{S}'(\beta) = -2\mathbf{Y}^t\mathbf{X} + 2\beta^t(\mathbf{X}^t\mathbf{X}) = \mathbf{0}. \quad (12)$$

3.4 Estimadores de la Regresión

- ▶ Como $\mathbf{X} \in \mathbf{M}_{n \times k+1}(\mathbb{R})$ es de rango completo:

$$\hat{\beta} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y}. \quad (13)$$

- ▶ El valor estimado $\hat{\mathbf{Y}}$, se calcula:

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\beta} = \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y} \in \mathbb{R}. \quad (14)$$

- ▶ Usando (14) podemos obtener la estimación de los errores $\hat{\epsilon}$:

$$\hat{\epsilon} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y} = (\mathbf{I} - \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t) \mathbf{Y}. \quad (15)$$

3.4 Estimadores de la Regresión

- ▶ Además de (15):

$$\mathbf{X}^t \hat{\epsilon} = \mathbf{X}^t (\mathbf{I} - \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t) \mathbf{Y} = \mathbf{0}, \quad (16)$$

- ▶ y :

$$\hat{\mathbf{Y}}^t \hat{\epsilon} = \hat{\beta} \mathbf{X}^t \hat{\epsilon} = \mathbf{0}. \quad (17)$$

- ▶ La variación cuadrática del error $\hat{\epsilon}$:

$$\hat{\epsilon}^t \hat{\epsilon} = \mathbf{Y}^t (\mathbf{I} - \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t) \mathbf{Y} = \mathbf{Y}^t \mathbf{Y} - \mathbf{Y}^t \mathbf{X} \hat{\beta}. \quad (18)$$

3.5 Coeficiente de Determinación

- De (18), tenemos:

$$\mathbf{Y}^t \mathbf{Y} = (\hat{\mathbf{Y}} + \hat{\epsilon})^t (\hat{\mathbf{Y}} + \hat{\epsilon}) = \hat{\mathbf{Y}}^t \hat{\mathbf{Y}} + \hat{\epsilon}^t \hat{\epsilon}, \quad (19)$$

- y de las varianzas:

$$\sum_{j=1}^n (\mathbf{Y}_j - \bar{\mathbf{Y}})^2 = \sum_{j=1}^n (\hat{\mathbf{Y}}_j - \bar{\mathbf{Y}})^2 + \sum_{j=1}^n \hat{\epsilon}_j^2. \quad (20)$$

- El coeficiente de determinación:

$$R^2 = 1 - \frac{\sum_{j=1}^n \hat{\epsilon}_j^2}{\sum_{j=1}^n (\mathbf{Y}_j - \bar{\mathbf{Y}})^2} = \frac{\sum_{j=1}^n (\hat{\mathbf{Y}}_j - \bar{\mathbf{Y}})^2}{\sum_{j=1}^n (\mathbf{Y}_j - \bar{\mathbf{Y}})^2}. \quad (21)$$

3.5 Coeficiente de Determinación

- ▶ De (21), el estadístico para el test de Hipótesis Nula:

$$H_0 : \beta_1 = \dots = \beta_k = 0, \quad (22)$$

- ▶ Es una F de Fisher de k y $n - k - 1$ grados de libertad:

$$F = \frac{\frac{R^2}{k}}{\frac{1-R^2}{n-k-1}}. \quad (23)$$

- ▶ Rechazo Hipótesis Nula si, con α el nivel de significación:

$$\mathbb{P}(F \geq f) \leq \alpha. \quad (24)$$

3.5 Coeficiente de Determinación

- ▶ Si $R^2 = 0$, significa que las variables no influyen en el modelo.
- ▶ Si $\hat{\beta} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y}$ entonces tenemos:
 - ▶ $E(\hat{\beta}) = \beta$, y
 - ▶ $Cov(\hat{\beta}) = \sigma^2 (\mathbf{X}^t \mathbf{X})^{-1}$.
- ▶ Para $\hat{\epsilon}$ tenemos:
 - ▶ $E(\hat{\epsilon}) = \mathbf{0}$, y
 - ▶ $Cov(\hat{\epsilon}) = \sigma^2 [\mathbf{I} - \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t]$.
- ▶ Luego de $E(\hat{\epsilon}^t \hat{\epsilon}) = (n - k - 1)\sigma^2$ tenemos:

$$s^2 = \frac{\hat{\epsilon}^t \hat{\epsilon}}{n - k - 1} = \frac{\mathbf{Y}^t [\mathbf{I} - \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t] \mathbf{Y}}{n - k - 1}. \quad (25)$$

3.5 Coeficiente de Determinación

- ▶ Para ver cuanto se ajusta el modelo, se puede ver el gráfico de distribución de los errores.
- ▶ Para ver el test de Hipótesis sobre los estimadores:

$$\begin{array}{ll} H_0 & \beta_j = \beta_j^* \\ H_1 & \beta_j \neq \beta_j^* \end{array}, \quad (26)$$

donde el valor β_j^* , es el valor particular que deseamos observar.

- ▶ Para ello observamos la variable aleatoria U_j , que tiene una distribución t_{n-k-1} :

$$U_j = \frac{\hat{\beta}_j - \beta_j^*}{\widehat{\sigma}_{\hat{\beta}_j}} = \frac{\hat{\beta}_j - \beta_j^*}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_j)}} = \frac{\hat{\beta}_j - \beta_j^*}{\sqrt{s^2((\mathbf{X}^t \mathbf{X})^{-1})_{jj}}}. \quad (27)$$

3.5 Coeficiente de Determinación

- ▶ Rechazamos la hipótesis nula cuando $|U_j| > c$, con c elegido desde un nivel de significación α_0 .
- ▶ Para un valor u_j observado, el valor de la cola esta dado por:

$$P(U_j > |u_j|) + P(U_j < -|u_j|) = p_{value}, \quad (28)$$

por lo tanto, rechazamos la hipótesis nula cuando $p_{value} \leq \alpha_0$.

- ▶ Para obtener un intervalo de confianza usamos:

$$\hat{\beta}_j \pm \left(t_{n-k-1} \left(\frac{\alpha}{2} \right) \right) \sqrt{s^2((\mathbf{X}^t \mathbf{X})^{-1})_{jj}}. \quad (29)$$

3.5 Coeficiente de Determinación

- ▶ Intervalo de confianza para la varianza del error:

$$\left(\frac{(n - k - 1)s^2}{\chi_{n-k-1}^2(1 - \frac{\alpha}{2})}, \frac{(n - k - 1)s^2}{\chi_{n-k-1}^2(\frac{\alpha}{2})} \right). \quad (30)$$

3.5 Coeficiente de Determinación

- ▶ Para estimar un valor, dado un vector x_0 , incluye el valor 1 si el modelo incluye la constante, en caso contrario no.

$$\hat{y}_0 = x_0^t \hat{\beta}. \quad (31)$$

Tenemos el valor:

$$h_0 = x_0^t (X^t X)^{-1} x_0. \quad (32)$$

El intervalo de confianza es:

$$\hat{y}_0 \pm s \sqrt{1 + h_0} \cdot t_{n-k-1} \left(1 - \frac{\alpha}{2}\right). \quad (33)$$

3.5 Coeficiente de Determinación

- ▶ Estadístico de Dubin Watson para ver autocorrelación entre el error:

$$D = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=2}^n e_i^2}. \quad (34)$$

El valor debe ser cercano a 2, para saber si es simétrico, si está cerca de 4 la autocorrelación es negativa, y si está cerca de 0, está correlacionado positivamente.

3.5 Coeficiente de Determinación

- ▶ Test de Breusch-Pagan :

$$H_0 : \text{Los errores tienen varianza constante.} \quad (35)$$

El valor del test es mayor que 0.05, entonces podemos asumir que los datos tienen varianza constante.