

2.3 REGRESIÓN MÚLTIPLE	18
2.1.1 Estimación por MCO de la función de regresión.....	18
2.1.2 Coeficiente de determinación R^2 corregido.....	23
2.1.3 Formas funcionales cuadráticas	25
2.1.4 Términos de interacción	28
2.1.5 Regresiones con variables estandarizadas	29
2.4 MODELIZACIÓN	30

* MRLM :

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} + \epsilon_i$$

Parte determinista
Parte aleatoria

- $K \equiv$ nº de variables independientes
- $K+1 \equiv$ nº de parámetros del modelo
- $n \equiv$ nº de muestras

$$\begin{aligned}
 n=1 \quad Y_1 &= \beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \dots + \beta_K X_{K1} + \epsilon_1 \\
 n=2 \quad Y_2 &= \beta_0 + \beta_1 X_{12} + \beta_2 X_{22} + \dots + \beta_K X_{K2} + \epsilon_2 \\
 &\vdots \\
 n \quad Y_n &= \beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \dots + \beta_K X_{Kn} + \epsilon_n
 \end{aligned}$$

$$Y = X \cdot \beta + \epsilon$$

$$Y = X \cdot \beta + \epsilon$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{21} & \dots & x_{k1} \\ 1 & x_{12} & x_{22} & \dots & x_{k2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1m} & x_{2m} & \dots & x_{km} \end{bmatrix} \cdot \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_m \end{bmatrix}$$

 $n \times 1$
 $n \times (k+1)$
 $(k+1) \times 1$
 $n \times 1$
 $n \times 1$

Objetivo MCO : minimizar $\sum \hat{\epsilon}_i^2$

$$\frac{\partial \sum \hat{\epsilon}_i^2}{\partial \beta_j} = 0 \Rightarrow \hat{\beta} = (X'X)^{-1} X'Y = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}$$

Y el mínimo de esta función:

$$\min_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki})^2$$

se consigue derivando respecto a cada parámetro e igualando a cero y operando se llega a las **k+1 ecuaciones normales**:

$$\sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki}) = \sum_{i=1}^n \hat{\varepsilon}_i = 0$$

$$\sum_{i=1}^n X_{ji} (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki}) = \sum_{i=1}^n X_{ij} \hat{\varepsilon}_i = 0$$

A partir de las **ecuaciones normales** se despejan los estimadores de los coeficientes y tenemos:

← tienen β_0

$$\begin{bmatrix} n & \sum X_{1i} & \sum X_{2i} & \dots & \sum X_{ki} \\ \sum X_{1i} & \sum X_{1i}^2 & \sum X_{1i}X_{2i} & \dots & \sum X_{1i}X_{ki} \\ \sum X_{2i} & \sum X_{2i}X_{1i} & \sum X_{2i}^2 & \dots & \sum X_{2i}X_{ki} \\ \dots & \dots & \dots & \dots & \dots \\ \sum X_{ki} & \sum X_{ki}X_{1i} & \sum X_{ki}X_{2i} & \dots & \sum X_{ki}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \dots & \dots & \dots & \dots \\ X_{k1} & X_{k2} & \dots & X_{kn} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix}$$

$$(X'X) \hat{\beta} = X' y$$

$\hat{\beta}_{MCO} = (X'X)^{-1} X'y$

$$\hat{\beta}_{MCO} = (X'X)^{-1} X'y = \begin{pmatrix} n & \sum_{i=1}^n X_{1i} & \dots & \sum_{i=1}^n X_{ki} \\ \sum_{i=1}^n X_{1i} & \sum_{i=1}^n X_{1i}^2 & \dots & \sum_{i=1}^n X_{1i}X_{ki} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n X_{ki} & \sum_{i=1}^n X_{ki}X_{1i} & \dots & \sum_{i=1}^n X_{ki}^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n Y_i X_{1i} \\ \vdots \\ \sum_{i=1}^n Y_i X_{ki} \end{pmatrix} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{pmatrix}$$

La matriz $X'X$ tiene dimensiones $(k+1) \times (k+1)$

$$X'X = \begin{pmatrix} n & \sum_{i=1}^n X_{1i} & \dots & \sum_{i=1}^n X_{ki} \\ \sum_{i=1}^n X_{1i} & \sum_{i=1}^n X_{1i}^2 & \dots & \sum_{i=1}^n X_{1i} X_{ki} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n X_{ki} & \sum_{i=1}^n X_{ki} X_{1i} & \dots & \sum_{i=1}^n X_{ki}^2 \end{pmatrix}$$

Cuadrada
 $(k+1)(k+1)$

- Donde n es el tamaño muestral.
- Las expresiones: $\sum_{i=1}^n X_{1i}^2 \dots \sum_{i=1}^n X_{ki}^2$ de la diagonal principal son las sumas de los cuadrados de las variable explicativas.
- Las expresiones: $\sum_{i=1}^n X_{1i} X_{ki} \dots \sum_{i=1}^n X_{ki} X_{1i}$ son las sumas de los productos cruzados de las variables explicativas.

En general las características que vimos en el análisis de regresión lineal simple se pueden extender al múltiple:

- $\sum_{i=1}^n \hat{y}_i \hat{\epsilon}_i = 0$ La estimación de la variable regresada y los residuos no están correlados.

Lo que implica que la covarianza es nula $cov(\hat{Y}_i, \hat{\epsilon}_i) = 0$

- La variables independientes X y los residuos también están incorrelacionados.

$$cov(X_{1i}, \hat{\epsilon}_i) = 0$$

Con los supuestos del modelo clásico de regresión lineal estos estimadores de MCO son lineales e insesgados, y dentro de éstos son los de mínima varianza (MELI).

ELI0 } Insesgado
BLUE } lineal
 } Eficiente

2.1.2 Coeficiente de determinación R^2 corregido

Una característica del modelo de regresión múltiple es que a medida que aumentamos el número de regresores X_j el coeficiente de determinación R^2 necesariamente aumenta salvo que el coeficiente estimado sea *exactamente* nulo. Entonces el R^2 nunca disminuye al incorporar nuevos regresores. Así que un incremento del R^2 no significa necesariamente que añadir una nueva variable realmente haya mejorado la calidad del ajuste de nuestro modelo. En realidad incluso si la nueva variable incluida en el modelo mejora nuestro ajuste, sabemos que necesariamente el R^2 de la nueva regresión estará artificialmente «inflado» por el mero hecho de incorporar un nuevo regresor. Por este motivo se utiliza el R^2 **corregido**, que ajusta por el número de coeficientes estimados y cuya definición es:

$$\bar{R}^2 = 1 - \frac{\frac{SCR}{n-k-1}}{\frac{SCT}{n-1}} = 1 - \frac{\hat{\sigma}^2}{S_Y^2}$$

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n-1}{n-k-1}$$

< R^2

Donde:

- $\hat{\sigma}^2$ es el estimador ~~ingresado~~ ^{inseguro} de la verdadera varianza de los residuos.
- S_Y^2 es la varianza muestral de Y.

A tener en cuenta:

- $\bar{R}^2 < R^2$ El coeficiente de determinación corregido siempre será menor que el coeficiente de determinación. → criterio de selección de modelos restringidos
- Añadir un nuevo regresor tiene dos efectos opuestos; el resultado final sobre \bar{R}^2 , el resultado final dependerá de cuál de los dos efectos es mayor:
 - Al disminuir la SCR → se incrementa la \bar{R}^2
 - El cociente $(n-1)/(n-k-1)$ aumenta → disminuir \bar{R}^2
- \bar{R}^2 puede ser negativo si los regresores en conjunto reducen la SCR en una cantidad tan pequeña que dicha reducción no logre superar el efecto del factor $(n-1)/(n-k-1)$.

Ejemplo 7: Consumo de las familias catalanas dedicadas a la hostelería y el turismo

Con datos de Cataluña y del sector de la hostelería, nos proponemos analizar el consumo de las familias catalanas.

Modelo poblacional: $\ln \text{consumo} = \beta_0 + \beta_1(\ln \text{ingresos}) + \varepsilon$

Y su estimación (FRM): $\widehat{\ln \text{consumo}} = 3,89 + 0,615(\ln \text{ingresos})$

$n = 95, R^2 = 0,3292.$

$$\begin{cases} \hat{\beta}_0 = 3,89 \\ \hat{\beta}_1 = 0,615 = \text{Eingresos} \end{cases}$$

$\frac{1}{x=0}$

Comparar \rightarrow ~~AX~~
 $\hookrightarrow R^2$

Este modelo es una regresión lineal simple. Es un modelo log-log, por tanto los coeficientes se interpretan directamente como las elasticidades. Un incremento del 1 % en los ingresos provoca que el consumo se incremente un 0,615 %.

Ampliando el modelo, ya que es muy probable que el consumo también tenga relación con el número de miembros en la familia (tamaño), ya que se espera que a medida que el tamaño aumenta, el consumo también aumenta, la nueva estimación:

$\widehat{\ln \text{consumo}}_i = 5,15 + 0,443 \cdot (\ln \text{ingreso}_i) + 0,1420 \cdot \text{tamaño}_i,$

$n = 95, R^2 = 0,4149.$

$$\begin{cases} \hat{\beta}_0 = 5,15 \\ \hat{\beta}_1 = 0,443 \\ \hat{\beta}_2 = 0,142 \end{cases}$$

ingresos \rightarrow log-log

tamaño \rightarrow log-lin

Al introducir la nueva variable, observamos que los coeficientes han cambiado, como hemos explicado anteriormente. Así que el término independiente y la elasticidad la variable ingresos cambia.

La nueva estimación nos aporta información sobre cómo influye el incremento, o decremento, del número de miembros en el consumo de las familias, dado un nivel determinado de ingresos.

Respecto a la variable tamaño el modelo es un modelo log-lineal así que la interpretación de los coeficientes estimados nos indica que, manteniendo constante el nivel de ingresos, es decir controlando el efecto del ingreso en el consumo, entonces el incremento de un miembro en la familia se prevé un incremento medio del 14,20 % del consumo familiar ($100 \times 0,1420 = 14,20$).

Por otro lado, el incremento de los ingresos en un 1 %, dado un tamaño familiar determinado, solo produce un incremento del 0,443 % del consumo, que contrasta con el 0,615 % de la expresión del modelo anterior de regresión lineal simple. Por tanto la introducción de nuevas variables (tamaño) afecta al resto de coeficientes de las variables del modelo (Ingreso).

Comparando los dos modelos podemos ver que el coeficiente de determinación del segundo modelo $R^2 = 0,4149$ es mayor que el del primero $R^2 = 0,3239$. Esto no podía ser de otro modo ya que como hemos comentado anteriormente el coeficiente de determinación aumenta al añadir regresores al modelo. Así que en este caso, no tiene sentido comparar estos coeficientes para decidir si la nueva variable añadida (tamaño) ha mejorado el ajuste del modelo, para poder determinar si el nuevo ajuste es mejor, necesitaríamos conocer el coeficiente de determinación corregido.

































