Modelo de regresión lineal múltiple

1 Ajuste mínimo-cuadrático del hiperplano de regresión

En el modelo de regresión múltiple que vamos a presentar se considera que el regresando¹ es una función lineal de k regresores (de los cuales k-1 corresponden a variables explicativas o a transformaciones de las mismas y uno corresponde al término independiente) y de una perturbación. Designando por Y_t al regresando, por X_{2t} , X_{3t} , ..., X_{kt} a los regresores y por u_t a la perturbación aleatoria, el modelo teórico de regresión lineal viene dado, para la observación genérica t-ésima, por la siguiente expresión:

Modelo de regresión múltiple

$$Y_{t} = \beta_{1} + \beta_{2} X_{2t} + \dots + \beta_{k} X_{kt} + u_{t} \qquad t = 1, 2, \dots, T$$
 (1)

Siendo T el tamaño de la muestra y dando valores a t desde t=1 hasta t=T, se obtiene el siguiente sistema de ecuaciones:

$$Y_{1} = \beta_{1} + \beta_{2}X_{21} + \beta_{3}X_{31} + \dots + \beta_{k}X_{k1} + u_{1}$$

$$Y_{2} = \beta_{1} + \beta_{2}X_{22} + \beta_{3}X_{32} + \dots + \beta_{k}X_{k2} + u_{2}$$

$$\dots \qquad \dots$$

$$Y_{T} = \beta_{1} + \beta_{2}X_{2T} + \beta_{3}X_{3T} + \dots + \beta_{k}X_{kT} + u_{T}$$

$$(2)$$

El sistema de ecuaciones anteriores se puede expresar de forma más compacta utilizando notación matricial. Así, vamos a denominar

$$\mathbf{y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_T \end{bmatrix} \qquad \mathbf{X} = \begin{bmatrix} 1 & X_{21} & X_{31} & \dots & X_{k1} \\ 1 & X_{22} & X_{32} & \dots & X_{k2} \\ \dots & \dots & \dots & \dots \\ 1 & X_{2T} & X_{3T} & \dots & X_{kT} \end{bmatrix} \qquad \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \dots \\ \boldsymbol{\beta}_k \end{bmatrix} \qquad \mathbf{u} = \begin{bmatrix} \boldsymbol{u}_1 \\ \boldsymbol{u}_2 \\ \dots \\ \boldsymbol{u}_T \end{bmatrix}$$

El modelo de regresión lineal múltiple (1) expresado en notación matricial es el siguiente:

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_T \end{bmatrix} = \begin{bmatrix} 1 & X_{21} & X_{31} & \dots & X_{k1} \\ 1 & X_{22} & X_{32} & \dots & X_{k2} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & X_{2T} & X_{3T} & \dots & X_{kT} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_k \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \\ \dots \\ u_T \end{bmatrix}$$
(3)

¹ El regresando puede ser o bien la variable endógena directamente o bien una transformación de la variable endógena (por ejemplo, el logaritmo de la variable endógena).

Si tenemos en cuenta las denominaciones dadas a vectores y matrices, el modelo de regresión lineal múltiple se puede expresar de forma compacta de la siguiente forma:

$$y = X\beta + u \tag{4}$$

donde, de acuerdo con la notación expuesta, y es un vector $T \times 1$, X es una matriz $T \times k$, β es un vector $k \times 1$ y u es un vector $T \times 1$.

El correspondiente modelo ajustado será el siguiente

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \tag{5}$$

El vector de residuos es igual a la diferencia entre valores observados y ajustados, es decir,

$$\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \tag{6}$$

Denominando S a la suma de los cuadrados de los residuos, se tiene que:

$$S = \hat{\mathbf{u}}'\hat{\mathbf{u}} = \begin{bmatrix} \hat{u}_1 & \hat{u}_2 & \dots & \hat{u}_T \end{bmatrix} \begin{bmatrix} \hat{u}_1 \\ \hat{u}_2 \\ \dots \\ \hat{u}_T \end{bmatrix} = \sum_{t=1}^T \hat{u}_t^2$$
 (7)

Teniendo en cuenta (6), se obtiene

$$S = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) =$$

$$= \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} - \mathbf{y}'\mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}}$$

$$= \mathbf{y}'\mathbf{y} - 2\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}}$$
(8)

Para llegar a la última igualdad de la expresión anterior se ha tenido en cuenta la igualdad entre los dos escalares siguientes

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

ya que un escalar es igual a su traspuesto, es decir,

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

Aplicar el criterio mínimo-cuadrático expuesto en el tema de regresión lineal simple es equivalente a minimizar el escalar S. Para ello se calcula la primera derivada de S con respecto al vector de coeficientes mínimo-cuadráticos, $\hat{\beta}$, en la expresión (8) y se iguala a 0^2 :

$$\frac{\partial S}{\partial \hat{\mathbf{\beta}}} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\mathbf{\beta}} = \mathbf{0}$$
 (9)

-

² Para la derivación de escalares, expresados mediante productos matriciales, respecto a un vector, véase el anexo 1 de *Econometría aplicada*.

Por tanto,

$$\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y} \tag{10}$$

Al sistema anterior se le denomina genéricamente sistema de ecuaciones normales del hiperplano. Cuando k=2 se obtiene el sistema de ecuaciones normales de la recta; cuando k=3 se obtiene el sistema de ecuaciones normales del plano; finalmente, cuando k>3 se obtiene específicamente el sistema de ecuaciones normales del hiperplano, el cuál no es susceptible de ser representado físicamente.

En notación matricial expandida, el sistema de ecuaciones normales es el siguiente:

$$\begin{bmatrix} T & \sum_{t=1}^{T} X_{2t} & \dots & \sum_{t=1}^{T} X_{kt} \\ \sum_{t=1}^{T} X_{2t} & \sum_{t=1}^{T} X_{2t}^{2} & \dots & \sum_{t=1}^{T} X_{2t} X_{kt} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{t=1}^{T} X_{kt} & \sum_{t=1}^{T} X_{kt} X_{2t} & \dots & \sum_{t=1}^{T} X_{kt}^{2} \end{bmatrix} \begin{bmatrix} \hat{\beta}_{1} \\ \hat{\beta}_{2} \\ \vdots \\ \hat{\beta}_{k} \end{bmatrix} = \begin{bmatrix} \sum_{t=1}^{T} Y_{t} \\ \sum_{t=1}^{T} X_{2t} Y_{t} \\ \vdots \\ \sum_{t=1}^{T} X_{kt} Y_{t} \end{bmatrix}$$
(11)

Obsérvese que:

- a) $\mathbf{X}'\mathbf{X}/T$ es la matriz de momentos muestrales de segundo orden, con respecto al origen, de los regresores.
- b). $\mathbf{X}'\mathbf{y}/T$ es el vector de momentos muestrales de segundo orden, con respecto al origen, entre el regresando y los regresores.

Para poder resolver el sistema (10) respecto a $\hat{\beta}$ unívocamente, se debe cumplir que el rango de la matriz $\mathbf{X}'\mathbf{X}$ sea igual a k. Si se cumple esta condición, se pueden premultiplicar ambos miembros de (10) por $\left[\mathbf{X}'\mathbf{X}\right]^{-1}$

$$[\mathbf{X}'\mathbf{X}]^{-1}\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = [\mathbf{X}'\mathbf{X}]^{-1}\mathbf{X}'\mathbf{y}$$

con lo cual se obtiene la expresión del vector de estimadores mínimo-cuadráticos:

$$\hat{\boldsymbol{\beta}} = \left[\mathbf{X}' \mathbf{X} \right]^{-1} \mathbf{X}' \mathbf{y} \tag{12}$$

S presenta un mínimo en $\hat{\beta}$, ya que la matriz de segundas derivadas, 2X'X, es definida positiva. Para comprobarlo, consideremos el vector \mathbf{a} , de orden $k \times 1$, distinto de cero. Entonces,

$$a'X'Xa = [Xa]'[Xa]$$

será el producto escalar del vector [Xa]' por su transpuesto, producto que no será negativo, por ser una suma de cuadrados. Si el rango de X'X es k, entonces queda garantizado que dicho producto es positivo. Por otra parte, al ser X'X definida positiva, se sigue, evidentemente, que 2X'X también es definida positiva.

2. Propiedades descriptivas en la regresión lineal múltiple

Las propiedades que se exponen a continuación son propiedades derivadas exclusivamente de la aplicación del método de estimación por mínimos cuadrados al modelo de regresión (1) en el que se incluye como primer regresor el término independiente.

1. La suma de los residuos mínimo-cuadráticos es igual a cero:

$$\sum_{t=1}^{T} \hat{u}_t = 0 \tag{13}$$

Demostración.

Por definición de residuo

$$\hat{u}_{t} = Y_{t} - \hat{Y}_{t} = Y_{t} - \hat{\beta}_{1} - \hat{\beta}_{2} X_{2t} - \dots - \hat{\beta}_{t} X_{t} \qquad t = 1, 2, \dots, T$$
 (14)

Si sumamos para las *T* observaciones se obtiene:

$$\sum_{t=1}^{T} \hat{u}_{t} = \sum_{t=1}^{T} Y_{t} - T \hat{\beta}_{1} - \hat{\beta}_{2} \sum_{t=1}^{T} X_{2t} - \dots - \hat{\beta}_{k} \sum_{t=1}^{T} X_{kt}$$
 (15)

Por otra parte, la primera ecuación del sistema de ecuaciones normales (11) es igual a

$$T\hat{\beta}_1 + \hat{\beta}_2 \sum_{t=1}^T X_{2t} + \dots + \hat{\beta}_k \sum_{t=1}^T X_{kt} = \sum_{t=1}^T Y_t$$
 (16)

Al comparar (15) y (16), se concluye que necesariamente debe cumplirse (13).

Obsérvese que, al cumplirse (13), se cumplirá también que

$$\sum_{t=1}^{T} Y_t = \sum_{t=1}^{T} \hat{Y}_t$$

y, al dividir por T, se obtiene

$$\overline{Y} = \overline{\hat{Y}} \tag{17}$$

2. El hiperplano de regresión pasa necesariamente por el punto $(\overline{Y}, \overline{X}_2, \cdots, \overline{X}_k)$.

Demostración.

En efecto, dividiendo la ecuación (16) por T se obtiene:

$$\overline{Y} = \hat{\beta}_1 + \hat{\beta}_2 \overline{X}_2 + \dots + \hat{\beta}_k \overline{X}_k \tag{18}$$

3. Los momentos de segundo orden entre cada regresor y los residuos son iguales a 0.

Para el conjunto de los regresores se puede expresar así:

$$\mathbf{X}'\hat{\mathbf{u}} = \mathbf{0} \tag{19}$$

En efecto,

$$X'\hat{u} = X' \Big\lceil y - X\hat{\beta} \Big\rceil = X'y - X'X\hat{\beta} = X'y - X'y = 0$$

Para llegar a la última igualdad se ha tenido en cuenta (6) y (10).

4. Los momentos de segundo orden entre $\hat{\mathbf{y}}$ y los residuos son 0, es decir,

$$\hat{\mathbf{y}}'\hat{\mathbf{u}} = \mathbf{0} \tag{20}$$

Demostración.

En efecto, si se tiene en cuenta (5) y (19) resulta que

$$\hat{\mathbf{y}}'\hat{\mathbf{u}} = \left[\mathbf{X}\hat{\boldsymbol{\beta}} \right]'\hat{\mathbf{u}} = \hat{\boldsymbol{\beta}}'\mathbf{X}'\hat{\mathbf{u}} = \hat{\boldsymbol{\beta}}'\mathbf{0} = \mathbf{0}$$

3 Hipótesis estadísticas básicas del modelo

I Hipótesis sobre la forma funcional

La relación entre el regresando, los regresores y la perturbación aleatoria es lineal:

$$Y_{t} = \beta_{1} + \beta_{2} X_{2t} + \dots + \beta_{k} X_{kt} + u_{t}$$
 $t = 1, 2, \dots, T$

o, en forma matricial,

$$y = X\beta + u \tag{21}$$

Il Hipótesis sobre el vector de perturbaciones aleatorias

La perturbación aleatoria u_t es una variable aleatoria no observable.

a) La esperanza matemática del vector de perturbaciones aleatorias es cero.

$$E(\mathbf{u}) = \mathbf{0} \tag{22}$$

b) Todas las perturbaciones aleatorias tienen la misma varianza, es decir, las perturbaciones son homoscedásticas

$$E(u_t)^2 = \sigma^2$$
 $t = 1, 2, ..., T$ (23)

c) Las perturbaciones aleatorias con distintos subíndices no están correlacionadas entre sí, es decir, las perturbaciones no están autocorrelacionadas:

$$E(u_t u_s) = 0 t \neq s (24)$$

La formulación de las hipótesis b) y c) permiten especificar la matriz de covarianzas del vector de perturbaciones. (La matriz de covarianzas de un vector que contiene T variables aleatorias es una matriz cuadrada y simétrica de orden $T \times T$, en cuya diagonal principal aparecen las varianzas de cada uno de los elementos del vector y fuera de la diagonal principal aparecen las covarianzas entre cada par de elementos.)

En concreto, la matriz de covarianzas del vector de perturbaciones es la siguiente:

$$E\left[\begin{bmatrix} \mathbf{u} - E(\mathbf{u}) \end{bmatrix} \begin{bmatrix} \mathbf{u} - E(\mathbf{u}) \end{bmatrix}' \right] = E\left[\begin{bmatrix} \mathbf{u} - \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} - \mathbf{0} \end{bmatrix}' \right] = E\left[\begin{bmatrix} \mathbf{u} \end{bmatrix} \begin{bmatrix} \mathbf{u} \end{bmatrix}' \right]$$

$$= E\left[\begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_T \end{bmatrix} \begin{bmatrix} u_1 & u_2 & \cdots & u_T \end{bmatrix} \right] = E\left[\begin{bmatrix} u_1^2 & u_1 u_2 & \cdots & u_1 u_T \\ u_2 u_1 & u_2^2 & \cdots & u_2 u_T \\ \vdots & \vdots & \ddots & \vdots \\ u_T u_1 & u_T u_2 & \cdots & u_T^2 \end{bmatrix}$$

$$= \begin{bmatrix} E(u_1^2) & E(u_1 u_2) & \cdots & E(u_1 u_T) \\ E(u_2 u_1) & E(u_2^2) & \cdots & E(u_2 u_T) \\ \vdots & \vdots & \ddots & \vdots \\ E(u_T u_1) & E(u_T u_2) & \cdots & E(u_T^2) \end{bmatrix} = \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix}$$

Para llegar a la última igualdad se ha tenido en cuenta que las varianzas de cada uno de los elementos del vector es constante e igual a σ^2 de acuerdo con (23) y que las covarianzas entre cada par de elementos es 0 de acuerdo con (24).

El resultado anterior se puede expresar de forma sintética:

$$E(\mathbf{u}\mathbf{u}') = \sigma^2 \mathbf{I} \tag{25}$$

A la anterior matriz se le denomina *matriz escalar*, ya que es igual a un escalar (σ^2 , en este caso) multiplicada por la matriz identidad.

d) La perturbación aleatoria tiene una distribución normal multivariante

Todas las hipótesis sobre el vector de perturbaciones se pueden formular de la siguiente forma:

$$\mathbf{u} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}) \tag{26}$$

III Hipótesis sobre el regresor X

a) La matriz de regresores, X, es una matriz fija.

De acuerdo con esta hipótesis, los distintos regresores del modelo toman los mismos valores para diversas muestras del regresando. Éste es un supuesto fuerte en el caso de las ciencias sociales, en el que es poco viable experimentar. Los datos se obtienen por observación, y no por experimentación. Para que dicho supuesto se cumpliera, los regresores deberían ser susceptibles de ser controlados por parte del investigador. Es importante señalar que los resultados que se obtienen utilizando este supuesto se mantendrían prácticamente idénticos si supusiéramos que los regresores son estocásticos, siempre que introdujéramos el supuesto adicional de independencia entre los regresores y la perturbación aleatoria. Este supuesto alternativo se puede formular así:

 a^*) La matriz de regresores, X, se distribuye independientemente del vector de perturbaciones aleatorias

$$E(\mathbf{X}'\mathbf{u}) = \mathbf{0} \tag{27}$$

b) La matriz de regresores, X, tiene rango k.

$$\rho(\mathbf{X}) = k \tag{28}$$

Recordemos que la matriz de regresores contiene k columnas, correspondientes a los k regresores del modelo, y T filas, correspondientes al número de observaciones. La hipótesis b) tiene dos implicaciones:

- 1. El número de observaciones, T, debe ser igual o mayor que el numero de regresores, k.
- 2. Todas las columnas de la matriz de regresores deben ser linealmente independientes, lo cual implica que no puede existir una relación lineal exacta entre ningún subconjunto de regresores. En caso contrario, el rango de la matriz **X** sería menor que *k*, y, por tanto, la matriz **X'X** sería singular (carecería de inversa), con lo cual no se podría determinar los valores del vector de estimadores de los parámetros del modelo. Si se diera este caso se dice que existe *multicolinealidad perfecta*. Si existe una relación lineal aproximada es decir, no exacta –, entonces se pueden obtener estimaciones de los parámetros, si bien la fiabilidad de los mismos quedaría afectada. En este último caso se dice que existe *multicolinealidad no perfecta*.

c) La matriz de regresores, X, no contiene errores de observación o de medida

Ésta es una hipótesis que raramente se cumple en la práctica, ya que los instrumentos de medición en economía son escasamente fiables (piénsese en la multitud de errores que es posible cometer en una recogida de información, mediante encuesta, sobre los presupuestos familiares). Aunque es difícil encontrar instrumentos para contrastar esta hipótesis, la naturaleza del problema y, sobre todo, la procedencia de los datos utilizados pueden ofrecer evidencia favorable o desfavorable a la hipótesis enunciada.

IV Hipótesis sobre el vector de parámetrosβ

El vector de parámetros β es constante.

Si no se adopta esta hipótesis el modelo de regresión sería muy complicado de manejar. En todo caso, puede ser aceptable que los parámetros del modelo se mantienen estables en el tiempo (si no se trata de períodos muy extendidos) o en el espacio (si está relativamente acotado).

4 Propiedades probabilísticas del modelo

Distribución del regresando

El regresando es función lineal del vector de perturbaciones aleatorias, que, por la hipótesis II d), sigue una distribución normal. Por lo tanto, el regresando, \mathbf{y} , seguirá también una distribución normal.

La esperanza matemática de \mathbf{y} , teniendo en cuenta la hipótesis II a)³, viene dada por

$$E(\mathbf{y}) = E[\mathbf{X}\boldsymbol{\beta} + \mathbf{u}] = \mathbf{X}\boldsymbol{\beta} + E(\mathbf{u}) = \mathbf{X}\boldsymbol{\beta}$$
 (29)

La matriz de varianzas covarianzas, teniendo en cuenta las hipótesis II a) a II c), serán

$$Var(\mathbf{y}) = E \left[(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' \right] = E(\mathbf{u}\mathbf{u}') = \sigma^2 \mathbf{I}$$
 (30)

En consecuencia, e1 regresando, y, tiene una distribución normal multivariante con vector de medias β y con una matriz de varianzas-covarianzas, $\sigma^2 \mathbf{I}$, escalar.

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}) \tag{31}$$

UU

³ Las hipótesis de los bloques III y IV se tendrán implícitamente en cuenta, aunque no se mencionen.