

Aspectos teóricos avanzados del modelo de regresión

Gabriel V. Montes-Rojas

Regresores estocásticos

- Hasta ahora asumimos que teníamos **regresores no estocásticos**. Este supuesto simplifica las pruebas de insesgadez y de cálculo de varianza, dado que las X s se pueden considerar números fijos.
- En la práctica si estamos dispuestos a asumir que tenemos una muestra aleatoria, esta muestra debiera contener factores aleatorios en todas las variables, $\{y_i, \mathbf{x}_i\}_{i=1}^N$.
- Como vamos a ver todo modelo de regresión tiene el objetivo de estimar una **esperanza condicional**, $E(y|\mathbf{x}) = \beta_1 x_1 + \dots + \beta_K x_K$, donde \mathbf{x} es el conjunto de todas las K variables explicativas, incluyendo una constante.
- “Condicionar” en una variable aleatoria es hacerla “fija”. En teoría de la probabilidad se basa en la definición de probabilidad condicional.

Esperanzas condicionales

- Toda variable aleatoria y se puede **descomponer** en dos partes ortogonales entre sí:

$$y = E(y|\mathbf{x}) + u,$$

donde

- (i) $E(u|\mathbf{x}) = 0$,
- (ii) $E(h(\mathbf{x})u) = 0$ para cualquier función $h(\cdot)$.

Prueba: (i) Definamos $u \equiv y - E(y|\mathbf{x})$. Tomando esperanzas $E(u|\mathbf{x}) = E(y|\mathbf{x}) - E(y|\mathbf{x}) = 0$. (ii) Usando la ley de esperanzas iteradas $E(h(\mathbf{x})u) = E(E(h(\mathbf{x})u|\mathbf{x})) = E(h(\mathbf{x})E(u|\mathbf{x})) = 0$.

- Un resultado importante es que $E(u|\mathbf{x}) = 0$ implica que $E(u) = 0$. Esto es por la propiedad de esperanzas iteradas que dice que $E_u(u) = E_{\mathbf{x}}[E_u(u|\mathbf{x})]$, donde la primera esperanza es con respecto a u y la segunda a \mathbf{x} .
- También, que $E(u|\mathbf{x}) = 0$ implica que $E(\mathbf{x}u) = 0$.
- Entonces, $\text{cov}(\mathbf{x}, u) \equiv E(\mathbf{x}u) - E(\mathbf{x})E(u) = 0$.
- Es decir, el supuesto $E(u|\mathbf{x}) = 0$ implica que los errores u de un modelo de regresión no están correlacionados con las \mathbf{x} .

Esperanzas condicionales

- La esperanza condicional es la solución al problema de minimización del valor esperado de las desviaciones al cuadrado, o sea

$$E(y|\mathbf{x}) = \arg \min_{m(\mathbf{x})} E((y - m(\mathbf{x}))^2).$$

Prueba: $(y - m(\mathbf{x}))^2 = ((y - E(y|\mathbf{x})) + (E(y|\mathbf{x}) - m(\mathbf{x})))^2 = (y - E(y|\mathbf{x}))^2 + (E(y|\mathbf{x}) - m(\mathbf{x}))^2 + 2(y - E(y|\mathbf{x}))(E(y|\mathbf{x}) - m(\mathbf{x}))$. Notemos que el primer término no depende de $m(\mathbf{x})$, mientras que el tercero se puede escribir como $u(\mathbf{x})(E(y|\mathbf{x}) - m(\mathbf{x})) = u(\mathbf{x})h(\mathbf{x})$. Si tomamos la esperanza condicional del tercero tenemos 0 por (ii).

- Sin embargo, no sabemos la forma funcional de $E(y|\mathbf{x})$.

Esperanzas condicionales

- Para cualquier variable aleatoria y , tenemos la proyección poblacional sobre el espacio generado por las \mathbf{x} , $r(\mathbf{x}) = \mathbf{x}\boldsymbol{\beta}$ donde $\boldsymbol{\beta} = \arg \min_{\mathbf{b}} E((y - \mathbf{x}\mathbf{b})^2)$.
- Si la esperanza condicional es lineal, entonces $E(y|\mathbf{x}) = \mathbf{x}\boldsymbol{\beta}$.
- Cada vez que corremos una regresión estamos estimando $E(y|\mathbf{x}) = \mathbf{x}\mathbf{b}$ asumiendo que es lineal en los parámetros. Conviene entonces decir que estamos estimando una esperanza condicional y levantar el supuesto de regresores no estocásticos.

Teorema de Gauss-Markov

- **Supuesto 1: Lineal en parámetros** La variable dependiente y se relaciona con X por una función lineal, i.e. $y = \beta_1 x_1 + \dots + \beta_K x_K + u$.
- **Supuesto 2: Muestreo aleatorio** $\{(y_i, x_{1i}, x_{2i}, \dots, x_{Ki})\}$, $i = 1, 2, \dots, N$ es una muestra aleatoria del modelo del Supuesto 1.
- **Supuesto 3: Ausencia de colinearidad perfecta en X** Para esto necesitamos que $(X'X)$ sea no singular o $\text{rango}(X'X)^{-1} = K$. Condición necesaria y suficiente para esto es que no haya una relación exacta entre los regresores.
- **Supuesto 4: Media condicional cero** $E[u|\mathbf{x}] = 0$

MCO es insesgado $E[\hat{\beta}_j|\mathbf{x}] = \beta_j$, $j = 1, 2, \dots, K$ o $E[\hat{\beta}|\mathbf{x}] = \beta$ donde β es el vector de todos los parámetros.

Prueba:...

Teorema de Gauss-Markov

- **Supuesto 5: Homocedasticidad** $\text{Var}[u|\mathbf{x}] = \sigma^2$

Teorema Gauss-Markov: Bajo los Supuestos 1-5, los estimadores MCO $(\hat{\beta}_1, \dots, \hat{\beta}_K)$ son los mejores estimadores lineales de β_1, \dots, β_K . Note: MEJOR significa mínima varianza.

Nociones básicas de teoría asintótica

- Consideremos una muestra aleatoria de tamaño N $\{x_1, x_2, \dots, x_N\}$ de una variable aleatoria $x \sim (\mu, \sigma^2)$. Esto es lo mismo que decir que es una muestra aleatoria de variables iid (independientes e idénticamente distribuidas).
- **Ley de los grandes números:** $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \xrightarrow{p} E[x] = \mu$ cuando $N \rightarrow \infty$.
Decimos que \bar{x} es **consistente** para μ . “ \xrightarrow{p} ” significa **convergencia en probabilidad**, $\bar{x} \xrightarrow{p} \mu$. También usamos la notación $\bar{x} = \mu + o_p(1)$, donde $o_p(1)$ significa que es un término que converge en probabilidad a 0.
- **Teorema central del límite:** $\frac{\sqrt{N}}{N} \sum_{i=1}^N x_i \xrightarrow{d} Normal(\mu, \sigma^2)$ cuando $N \rightarrow \infty$. “ \xrightarrow{d} ” significa **convergencia en distribución**.
- Decimos que \bar{x} es asintóticamente normal: $\sqrt{N}(\bar{x} - \mu) \xrightarrow{d} Normal(0, \sigma^2)$ cuando $N \rightarrow \infty$.

Identificación

- Supuesto **OLS.0** (muestra aleatoria): El modelo poblacional es $y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_K x_{Ki} + u_i, i = 1, 2, \dots$ (también escrito como $y_i = \mathbf{x}_i \boldsymbol{\beta} + u_i$ con $\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ki})$) donde la $\{(y_i, \mathbf{x}_i, u_i) : i = 1, 2, \dots\}$ es una secuencia de vectores aleatorios iid.
- Supuesto **OLS.1** (condición de ortogonalidad poblacional): $E(\mathbf{x}'u) = 0$
 - Dado que \mathbf{x} tiene una constante, OLS.1 es equivalente a decir que $E(u) = 0$ y que $Cov(x_j, u) = 0, j = 1, 2, 3, \dots, K$.
 - Una condición suficiente para OLS.1 es $E(u|\mathbf{x}) = 0$:
 $E(\mathbf{x}'u) = E(\mathbf{x}'E(u|\mathbf{x})) = 0$.
- Supuesto **OLS.2**: $\text{rango } E(\mathbf{x}'\mathbf{x}) = K$
 - Dado que $E(\mathbf{x}'\mathbf{x})$ es una matrix simétrica $K \times K$, OLS.2 es equivalente a decir que $E(\mathbf{x}'\mathbf{x})$ es positiva definida y que $Cov(x_j, u) = 0, j = 1, 2, 3, \dots, K$.

Identificación

Bajo los supuestos OLS.0, OLS.1 y OLS.2, el vector de parámetros β está **identificado**. En este contexto significa que se puede expresar como momentos poblacionales de variables observadas.

$$\beta = [E(\mathbf{x}'\mathbf{x})]^{-1}E(\mathbf{x}'y).$$

Prueba: De OLS.0 reemplazar $y = \mathbf{x}\beta + u$ en $\mathbf{x}'y = \mathbf{x}'(\mathbf{x}\beta + u)$. Aplicar esperanzas a ambos lados de la igualdad, $E(\mathbf{x}'y) = E(\mathbf{x}'\mathbf{x}\beta + \mathbf{x}'u)$. Usando OLS.1 y OLS.2 llegamos al resultado. QED

Consistencia

Podemos reescribir el estimador MCO como

$$\begin{aligned}\hat{\beta} &= \left(N^{-1} \sum_i^N \mathbf{x}'_i \mathbf{x}_i \right)^{-1} \left(N^{-1} \sum_i^N \mathbf{x}'_i y_i \right) = \beta + \left(N^{-1} \sum_i^N \mathbf{x}'_i \mathbf{x}_i \right)^{-1} \left(N^{-1} \sum_i^N \mathbf{x}'_i u_i \right) \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}\end{aligned}$$

Bajo los supuestos OLS.0, OLS.1 y OLS.2, el estimador MCO $\hat{\beta}$ de una muestra aleatoria es **consistente** para β . Es decir, $\hat{\beta} \xrightarrow{P} \beta$, $N \rightarrow \infty$.

Prueba: Usando la **ley de los grandes números**, $\left(N^{-1} \sum_i^N \mathbf{x}'_i \mathbf{x}_i \right) \xrightarrow{P} E(\mathbf{x}'\mathbf{x})$ y $\left(N^{-1} \sum_i^N \mathbf{x}'_i y_i \right) \xrightarrow{P} E(\mathbf{x}'y)$. Usando identificación se prueba. [También se puede probar con $\left(N^{-1} \sum_i^N \mathbf{x}'_i u_i \right) \xrightarrow{P} E[\mathbf{x}'u] = \mathbf{0}$.] QED

Normalidad asintótica

La distribución asintótica de $\hat{\beta}$ depende de

$$\sqrt{N}(\hat{\beta} - \beta) = \left(N^{-1} \sum_i^N \mathbf{x}'_i \mathbf{x}_i \right)^{-1} \left(N^{-1/2} \sum_i^N \mathbf{x}'_i u_i \right).$$

- Probar que $\left(N^{-1} \sum_i^N \mathbf{x}'_i \mathbf{x}_i \right)^{-1} - \mathbf{A}^{-1} = o_p(1)$, donde $\mathbf{A} \equiv E(\mathbf{x}'\mathbf{x})$ es una matriz simétrica $K \times K$. Definamos $\hat{\mathbf{A}} \equiv N^{-1} \sum_i^N \mathbf{x}'_i \mathbf{x}_i$ como el estimador de esta matriz.
- $\{\mathbf{x}'_i u_i\} : i = 1, 2, \dots\}$ es una secuencia iid con media cero y varianza finita. Entonces por una aplicación del **teorema central del límite**,

$$N^{-1/2} \sum_i^N \mathbf{x}'_i u_i \xrightarrow{d} \text{Normal}(\mathbf{0}, \mathbf{B}), \quad N \rightarrow \infty$$

donde $\mathbf{B} \equiv \text{Var}(\mathbf{x}'u) = E(\mathbf{x}'uu'\mathbf{x}) - E(\mathbf{x}'u)E(u'\mathbf{x}) = E(\mathbf{x}'uu'\mathbf{x}) = E(u^2\mathbf{x}'\mathbf{x})$ es una matriz $K \times K$.

Normalidad asintótica

- Entonces,

$$\sqrt{N}(\hat{\beta} - \beta) = \mathbf{A}^{-1} \left(N^{-1/2} \sum_i^N \mathbf{x}'_i u_i \right) + o_p(1)$$

$o_p(1)$ es que se hace 0 asintóticamente en probabilidad.

- En base a los resultados anteriores tenemos

$$\sqrt{N}(\hat{\beta} - \beta) \xrightarrow{d} \text{Normal}(\mathbf{0}, \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1}),$$

donde estamos haciendo uso de la simetría de \mathbf{A} . Aparece acá la típica forma sandwich en la varianza de los estimadores.

Asumimos homoscedasticidad:

Supuesto **OLS.3**: $\text{Var}(\mathbf{x}'u) = E(u^2 \mathbf{x}'\mathbf{x}) = \sigma^2 \mathbf{A}$.

Bajo los supuestos OLS.0, OLS.1, OLS.2 y OLS.3, tenemos cuando $N \rightarrow \infty$,

$$\sqrt{N}(\hat{\beta} - \beta) \xrightarrow{d} \text{Normal}(\mathbf{0}, \sigma^2 \mathbf{A}^{-1})$$

como resultado del **teorema central del límite**. Es decir, los estimadores de MCO son **asintóticamente normales**.

Normalidad asintótica

- Podemos definir la varianza asintótica como $AVar(\hat{\beta}) = \sigma^2 \mathbf{A}^{-1}$. Sin embargo todavía tenemos que proponer un estimador de la varianza, $\widehat{AVar}(\hat{\beta}) = \hat{\sigma}^2 \hat{\mathbf{A}}^{-1}$.
- Lema: $\hat{\sigma}^2 ((\mathbf{X}'\mathbf{X})/N)^{-1} \xrightarrow{P} \sigma^2 \mathbf{A}^{-1}$ cuando $N \rightarrow \infty$, donde $\hat{\sigma}^2 \equiv N^{-1} \sum_i^N \hat{u}_i^2$.

Prueba: Primero, tenemos que probar que $\hat{\sigma}^2 \xrightarrow{P} \sigma^2$. Para ello consideremos $\hat{\sigma}^2 \equiv N^{-1} \sum_i^N \hat{u}_i^2 = N^{-1} (\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y})'(\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}) = N^{-1}(\mathbf{M}_X \mathbf{y})'(\mathbf{M}_X \mathbf{y})$, donde $\mathbf{M}_X = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ que es simétrica e idempotente. Ver **Proyecciones** para una definición formal. Ahora, $\mathbf{y} = \mathbf{X}\beta + \mathbf{u}$ y $\mathbf{M}_X \mathbf{X} = \mathbf{0}$. Entonces, $\hat{\sigma}^2 = N^{-1} \mathbf{u}'\mathbf{u} - N^{-1} \mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} = N^{-1} \sum_i^N u_i^2 + o_p(1) = \sigma^2 + o_p(1)$.

Por otro lado, como teníamos más arriba, $\hat{\mathbf{A}} \equiv N^{-1} \sum_i^N \mathbf{x}_i' \mathbf{x}_i$ es el estimador de \mathbf{A} . Tenemos que probar que es consistente. Para ello hacemos uso del supuesto OLS.0 (muestra aleatoria). Luego necesitamos que \mathbf{A} sea invertible, que lo obtenemos con OLS.2. *QED*

Teorema de Frisch-Waugh-Lovell

Consideremos los siguientes modelos de regresión:

$$A. y = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K + u$$

$$B. y = \gamma_2 x_2 + \gamma_3 x_3 + \dots + \gamma_K x_K + v$$

$$C. x_1 = \delta_2 x_2 + \delta_3 x_3 + \dots + \delta_K x_K + e$$

- Computar los residuos del modelo B (\hat{v}) y C (\hat{e}).
- Correr la siguiente regresión auxiliar: $\hat{v} = \alpha \hat{e} + \text{residuo}$
- Chequear que $\hat{\alpha}_1 = \hat{\beta}_1$
- ¿Cómo se interpreta este resultado?

Teorema de Frisch-Waugh-Lovell

En notación matricial, podemos escribir el modelo de regresión múltiple en forma general como

$$y = X\beta + u = X_1\beta_1 + X_2\beta_2 + u$$

donde $X = [X_1 \ X_2]$ y $\beta = [\beta_1' \ \beta_2']'$.

Ahora construyamos $M_2y = M_2X_1\beta_1 + M_2u$ donde M_2 es la proyección residual de X_2 , es decir, $M_2 = X_2(X_2'X_2)^{-1}X_2'$. Ver [Proyecciones](#) para una definición formal.

El teorema de FWL muestra que

$$\hat{\beta}_1 = (X_1'M_2X_1)^{-1}X_1'M_2y = (\tilde{X}_1'\tilde{X}_1)^{-1}\tilde{X}_1'\tilde{y},$$

donde $\tilde{X}_1 = M_2X_1$ e $\tilde{y} = M_2y$, es decir, las proyecciones residuales de X_1 e y con respecto a X_2 .

Teorema de Frisch-Waugh-Lovell

Prueba: Consideremos la siguiente igualdad

$$\mathbf{y} = \mathbf{P}_X \mathbf{y} + \mathbf{M}_X \mathbf{y} = \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \mathbf{M}_X \mathbf{y},$$

donde tenemos que $\mathbf{P}_X \mathbf{y} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2$ dado que $\hat{\boldsymbol{\beta}} = [\hat{\boldsymbol{\beta}}_1' \ \hat{\boldsymbol{\beta}}_2']'$.
Multipliquemos ambos lados de la igualdad por $\mathbf{X}'_1 \mathbf{M}_2$. Entonces obtenemos

$$\mathbf{X}'_1 \mathbf{M}_2 \mathbf{y} = \mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \mathbf{X}'_1 \mathbf{M}_2 \mathbf{M}_X \mathbf{y}.$$

Ahora usamos el resultado $\mathbf{M}_2 \mathbf{X}_2 = \mathbf{0}$ y $\mathbf{M}_X \mathbf{M}_2 \mathbf{X}_1 = \mathbf{0}$. Entonces,

$$\mathbf{X}'_1 \mathbf{M}_2 \mathbf{y} = \mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1.$$

Resolviendo por $\hat{\boldsymbol{\beta}}_1$,

$$\hat{\boldsymbol{\beta}}_1 = (\mathbf{X}'_1 \mathbf{M}'_2 \mathbf{M}_2 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{M}'_2 \mathbf{M}_2 \mathbf{y} = (\mathbf{X}'_1 \mathbf{M}'_2 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{M}_2 \mathbf{y},$$

donde usamos el resultado que \mathbf{M}_2 es una matriz simétrica e idempotente. *QED*

Algebra de MCO (nota para la slide anterior)

Proyección ortogonal: Definamos la matriz $N \times N$ $\mathbf{P}_X \equiv \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ como la matriz que proyecta y en el espacio generado por \mathbf{X} . Tenemos que $\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{P}_X\mathbf{y}$, los valores predichos.

Proyección residual: Definamos la matriz $N \times N$ $\mathbf{M}_X \equiv \mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ como la proyección de y en el complemento del espacio de \mathbf{X} . Notemos que $\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I}_N - \mathbf{P}_X)\mathbf{y} = \mathbf{M}_X\mathbf{y}$, residuos de la regresión.

Notar que $\mathbf{P}_X\mathbf{M}_X = \mathbf{0}$, es decir, son ortogonales. Esto significa que $\mathbf{y} = (\mathbf{P}_X + \mathbf{M}_X)\mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{u}} = \mathbf{P}_X\mathbf{y} + \mathbf{M}_X\mathbf{y}$, es decir, el vector \mathbf{y} se descompone en dos partes ortogonales entre sí.

Además ambas matrices son simétricas, $\mathbf{P}'_X = \mathbf{P}_X$, $\mathbf{M}'_X = \mathbf{M}_X$, e idempotentes, $\mathbf{P}_X\mathbf{P}_X = \mathbf{P}_X$, $\mathbf{M}_X\mathbf{M}_X = \mathbf{M}_X$.

El problema con la multicolinealidad

- La multicolinealidad perfecta invalida el modelo de regresión, ej., $X'X$ no es invertible.
- La multicolinealidad imperfecta no es necesariamente un problema. La mayoría de las variables de control están correlacionadas entre sí (ej., educación y experiencia), y eso es justamente lo que hace al modelo de regresión múltiple.
- El problema es cuando la correlación es muy alta (ej., dos índices de precios del mismo país). El resultado es que la significatividad individual es baja (test t), pero la significatividad global del modelo (test F, R^2) es alta. Eso se puede ver con el Teorema de FWL.