

Multicolinealidad y Variables Cualitativas

Fundamentos de Econometría

Juan Palomino¹

¹Magister en Economía Aplicada con Mención Estudios Regionales
juan.palominoh@pucp.pe

Departamento de Economía



Índice

1 Multicolinealidad

- Multicolinealidad Perfecta
- Multicolinealidad Imperfecta

2 Variables Cualitativas

- Variables Dummy
- Interacciones
- Variables Categóricas

Índice

1 Multicolinealidad

- Multicolinealidad Perfecta
- Multicolinealidad Imperfecta

2 Variables Cualitativas

- Variables Dummy
- Interacciones
- Variables Categóricas

Multicolinealidad perfecta

- **Definición:** una de las variables explicativas es una combinación lineal exacta de las otras explicativas.
- Por ejemplo, dado el modelo poblacional:

$$PBI_i = \beta_0 + \beta_1 export + \beta_2 import_i + \beta_3 balance + \varepsilon_i$$

- Ocurriría multicolinealidad perfecta si podemos expresar $balance = export - import$.
- En este caso β_1 no podrá medir el efecto de incrementar $balance$ manteniendo constante $export$ e $import$ dada su relación lineal.

Multicolinealidad perfecta

- No podemos encontrar de forma única $\hat{\beta}_{OLS} = (X'X)^{-1}X'Y$

$$\text{rango}(X) < K \Rightarrow |X'X| = 0 \Rightarrow \nexists (X'X)^{-1}$$

- ¿Cómo detectarlo?
 - Los programas se quejaron de que no podemos invertir la matriz $(X'X)$
- ¿Cómo corregirlo?
 - Se deben a errores del investigador al introducir las explicativas
 - Al aparecer mensaje de error, corregiremos las explicativas.

Multicolinealidad perfecta

- Corrección: en el ejemplo de balanza comercial sabemos:
 $balance = export - import$

$$PBI_i = \beta_0 + \beta_1 export + \beta_2 import_i + \beta_3 balance + \varepsilon_i$$

- Excluimos la variable *balance*:

$$PBI_i = \beta_0 + \beta_1 export + \beta_2 import_i + \varepsilon_i$$

Índice

1 Multicolinealidad

- Multicolinealidad Perfecta
- Multicolinealidad Imperfecta

2 Variables Cualitativas

- Variables Dummy
- Interacciones
- Variables Categóricas

Multicolinealidad Imperfecta

- **Definición:** la correlación entre las variables explicativas es alta, pero no perfecta.
- Por ejemplo:

$$\text{Consumo}_i = \beta_1 + \beta_2 \text{Ingreso}_i + \beta_3 \text{Riqueza}_i + \varepsilon_i$$

- No hay multicolinealidad perfecta, pero es probable que riqueza e ingresos estén altamente correlacionadas.

Multicolinealidad Imperfecta

- Podemos encontrar de forma única $\hat{\beta}_{OLS} = (X'X)^{-1}X'Y$

$$rango(X) = K \Rightarrow |X'X| \neq 0 \Rightarrow \exists (X'X)^{-1}$$

- Estimadores cumplen las propiedades de MCO.
- ¿Qué problemas genera en la estimación?
 - Para entenderlo, consideremos un modelo de regresión lineal con k variables:

$$Y_i = \beta_1 + \beta_2 W_{2i} + \dots + \beta_{k-1} W_{k-1i} + \beta_z Z_i + \varepsilon_i$$

- Y definimos:

$$R_z^2 = 1 - \frac{SCR}{STC} = 1 - \frac{\varepsilon_z' \varepsilon_z}{\sum_{i=1}^n (Z_i - \bar{Z})^2}$$

Multicolinealidad Imperfecta

- ¿Qué problemas genera en la estimación?
 - Podemos demostrar por teorema de Frisch-Waugh que:

$$\hat{\beta}_Z = (z' M_W z)^{-1} z' M_W y$$

- Entonces:

$$\text{Var}(\hat{\beta}_Z) = \sigma^2 (z' M_W z)^{-1}$$

$$= \frac{\sigma^2}{z' M_W' M_W z}$$

$$\text{Var}(\hat{\beta}_Z) = \frac{\sigma^2}{\hat{\varepsilon}_z' \hat{\varepsilon}_z}$$

Multicolinealidad Imperfecta

- Sabiendo que $\hat{\varepsilon}_Z$ son los residuos de la regresión de Z_i contra las variables W . Sabiendo $\varepsilon'_z \varepsilon_z$ tenemos que:

$$Var(\hat{\beta}_Z) = \frac{\sigma^2}{(1 - R_z^2) \sum_{i=1}^n (Z_i - \bar{Z})^2}$$

- Implicaciones:
 - Si $R_z^2 \rightarrow 1 \Rightarrow Var(\hat{\beta}_Z) \rightarrow \infty$
 - Si Z no se relacionara con las demás variables, entonces $R_z^2 = 0$.
 - Estimación imprecisa e intervalos de confianza muy grandes.

Multicolinealidad Imperfecta

- Pruebas estadísticas informales:
 - Análisis de la matriz de correlación
 - Analizar el R^2 , su test F asociado y las pruebas t individuales.
 - Si el R^2 es alto y el test F asociado es estadísticamente significativo pero las pruebas t individuales son bajas, muy probablemente exista multicolinealidad.

Multicolinealidad Imperfecta

- Pruebas estadísticas formales:
 - Analizar el VIF (variance inflation factor): mide la redundancia que una variable inserta en el modelo en términos de la información que comparte con las otras variables:

$$\begin{aligned}
 \text{Var}(\hat{\beta}_Z) &= \frac{\sigma^2}{(1 - R_z^2) \sum_{i=1}^n (Z_i - \bar{Z})^2} \\
 &= \frac{1}{(1 - R_z^2)} \frac{\sigma^2}{\sum_{i=1}^n (Z_i - \bar{Z})^2} \\
 &= \text{VIF} \frac{\sigma^2}{\sum_{i=1}^n (Z_i - \bar{Z})^2}
 \end{aligned}$$

Multicolinealidad Imperfecta

- Existen dos reglas para interpretar el VIF :
 - Si el VIF de una variable es mayor a 10, podemos afirmar que esta inserta al modelo una considerable inflación de la varianza.
 - Si el promedio de los VIF es considerablemente mayor a 1, podemos afirmar que el modelo sufre un problema de multicolinealidad.

Índice

1 Multicolinealidad

- Multicolinealidad Perfecta
- Multicolinealidad Imperfecta

2 Variables Cualitativas

- Variables Dummy
- Interacciones
- Variables Categóricas

Variables Cualitativas

- Tenemos las siguientes variables cualitativas:

$$Sexo = \left\{ \begin{array}{l} Hombre \\ Mujer \end{array} \right\}$$

- Para la variable sexo podemos tener:

$$S_{1i} = \left\{ \begin{array}{l} 1 \text{ si es hombre} \\ 0 \text{ si es mujer} \end{array} \right\} \text{ o } S_{2i} = \left\{ \begin{array}{l} 1 \text{ si es mujer} \\ 0 \text{ si es hombre} \end{array} \right\}$$

Variables Cualitativas

- Supongamos que el modelo es:

$$\ln(\text{ingreso})_i = \beta_0 + \beta_1 \text{educ}_i + \varepsilon_i$$

- Incorporar la variable sexo:

$$\ln(\text{ingreso})_i = \beta_0 + \beta_1 \text{educ}_i + \beta_2 S_{1i} + \beta_3 S_{2i} + \varepsilon_i$$

- Este modelo no puede ser estimado. La manera correcta es:

$$\ln(\text{ingreso})_i = \beta_0 + \beta_1 \text{educ}_i + \beta_2 S_{1i} + \varepsilon_i$$

- En la función de regresión poblacional

$$E[\text{ingreso}_i \mid S_{1i} = 1, \text{educ}] = (\beta_0 + \beta_2) + \beta_1 \text{educ}_i$$

$$E[\text{ingreso}_i \mid S_{1i} = 0, \text{educ}] = \beta_0 + \beta_1 \text{educ}_i$$

Variables Cualitativas

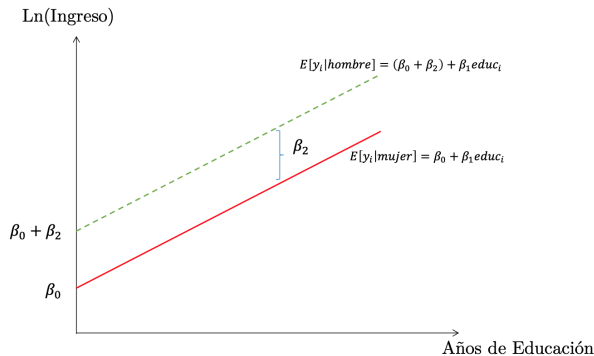


Figura: FRP de los hombres y mujeres con cambios en intercepto

Índice

1 Multicolinealidad

- Multicolinealidad Perfecta
- Multicolinealidad Imperfecta

2 Variables Cualitativas

- Variables Dummy
- **Interacciones**
- Variables Categóricas

Variables Cualitativas: Dummy Interactivas

- Supongamos que el modelo es:

$$\ln(\text{ingreso})_i = \beta_0 + \beta_1 \text{educ}_i + \beta_2 S_{1i} + \beta_3 (S_{1i} \times \text{educ})_i + \varepsilon_i$$

- El valor esperado condicional es:

$$E[\text{ingreso}_i \mid S_{1i} = 1, \text{educ}] = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \text{educ}_i$$

$$E[\text{ingreso}_i \mid S_{1i} = 0, \text{educ}] = \beta_0 + \beta_1 \text{educ}_i$$

Variables Cualitativas: Dummy Interactivas

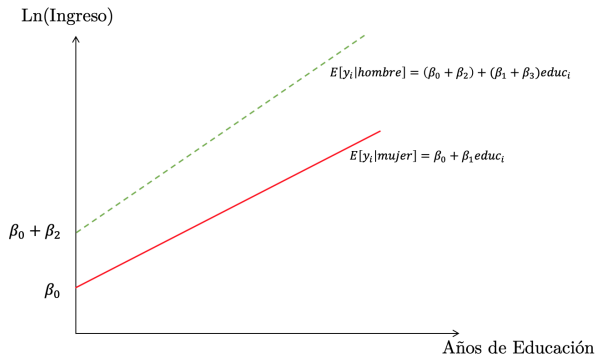


Figura: FRP de hombres y mujeres con cambios en intercepto y pendiente

Índice

1 Multicolinealidad

- Multicolinealidad Perfecta
- Multicolinealidad Imperfecta

2 Variables Cualitativas

- Variables Dummy
- Interacciones
- Variables Categóricas

Variables Cualitativas

- Tenemos las siguientes variables cualitativas:

$$Región = \begin{Bmatrix} Costa \\ Sierra \\ Selva \end{Bmatrix}$$

- Para la variable región podemos tener:

$$R_{1i} = \begin{Bmatrix} 1 \text{ si vive en la Costa} \\ 0 \text{ en otro caso} \end{Bmatrix}$$

$$R_{2i} = \begin{Bmatrix} 1 \text{ si vive en la Sierra} \\ 0 \text{ en otro caso} \end{Bmatrix}$$

$$R_{3i} = \begin{Bmatrix} 1 \text{ si vive en la Selva} \\ 0 \text{ en otro caso} \end{Bmatrix}$$

Variables Cualitativas

- Supongamos que el modelo es:

$$\ln(\text{ingreso})_i = \beta_0 + \beta_1 \text{educ}_i + \beta_2 S_{1i} + \beta_3 R_{2i} + \beta_4 R_{3i} + \varepsilon_i$$

Variables Cualitativas

- Mujeres de la Costa:

$$E[\text{ingreso}_i \mid S_{1i} = 0, R_{2i} = 0, R_{3i} = 0, \text{educ}] = \beta_0 + \beta_1 \text{educ}_i$$

- Hombres de la Costa:

$$E[\text{ingreso}_i \mid S_{1i} = 1, R_{2i} = 0, R_{3i} = 0, \text{educ}] = (\beta_0 + \beta_2) + \beta_1 \text{educ}_i$$

- Mujeres de la Sierra:

$$E[\text{ingreso}_i \mid S_{1i} = 0, R_{2i} = 1, R_{3i} = 0, \text{educ}] = (\beta_0 + \beta_3) + \beta_1 \text{educ}_i$$

- Hombres de la Sierra:

$$E[\text{ingreso}_i \mid S_{1i} = 1, R_{2i} = 1, R_{3i} = 0, \text{educ}] = (\beta_0 + \beta_2 + \beta_3) + \beta_1 \text{educ}_i$$

- Mujeres de la Selva:

$$E[\text{ingreso}_i \mid S_{1i} = 0, R_{2i} = 0, R_{3i} = 1, \text{educ}] = (\beta_0 + \beta_4) + \beta_1 \text{educ}_i$$

- Hombres de la Selva:

$$E[\text{ingreso}_i \mid S_{1i} = 1, R_{2i} = 0, R_{3i} = 1, \text{educ}] = (\beta_0 + \beta_2 + \beta_4) + \beta_1 \text{educ}_i$$

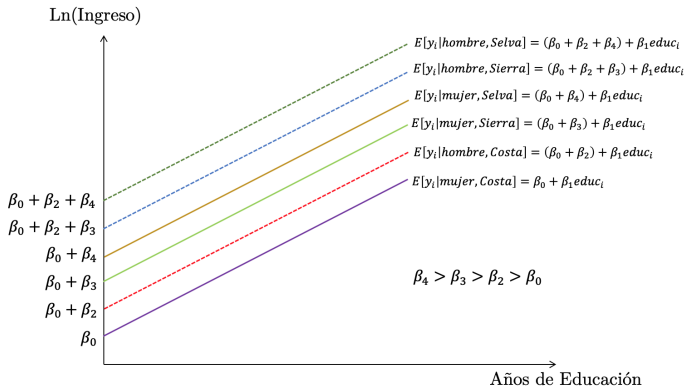


Figura: FRP de los hombres y mujeres en Regiones

Variables Cualitativas

- Supongamos que el modelo es:

$$\ln(\text{ingreso})_i = \beta_0 + \beta_1 \text{educ}_i + \beta_2 S_{1i} + \beta_3 R_{2i} + \beta_4 R_{3i} + \beta_5 (S_1 \times R_2)_i + \beta_6 (S_1 \times R_3)_i + \varepsilon_i$$

Variables Cualitativas

■ Diferencial Sierra-Costa (hombres)

$$\begin{aligned} E[y_i \mid S_{1i} = 1, R_{2i} = 1, R_{3i} = 0, educ] - E[y_i \mid S_{1i} = 1, R_{2i} = 0, R_{3i} = 0, educ] \\ (\beta_0 + \beta_1 educ_i + \beta_2 + \beta_3 + \beta_5) - (\beta_0 + \beta_1 educ_i + \beta_2) \\ = \beta_3 + \beta_5 \end{aligned}$$

■ Diferencial Sierra-Costa (mujeres)

$$\begin{aligned} E[y_i \mid S_{1i} = 0, R_{2i} = 1, R_{3i} = 0, educ] - E[y_i \mid S_{1i} = 0, R_{2i} = 0, R_{3i} = 0, educ] \\ (\beta_0 + \beta_1 educ_i + \beta_3) - (\beta_0 + \beta_1 educ_i) \\ = \beta_3 \end{aligned}$$