

# Selección en observables

En otras palabras: el término de error (el efecto no observable en  $Y$ ) no presenta diferencias entre el grupo tratado y no tratado  
El error es igual a si fue tratado o no  $\Rightarrow$  Medida condicional  $\rightarrow$  INDEPENDENCIA

$$E(Y_{1i} | T_i = 1) - E(Y_{0i} | T_i = 1) + \left[ E(Y_{1i} | T_i = 0) - E(Y_{0i} | T_i = 0) \right]$$

↑  
si tengo S.S.  
subestimar  
sobre estimar

## SESGO DE SELECCIÓN

Factores no observados no están relacionados con si se trató o no

### ALEATORIZACIÓN

en un experimento controlado

sin experimento

SELECCIÓN OBSERVABLES

Factores observados

$u$  correlacionado a  $X_i$

$Y_i$  depende de la relación de  $X_i$  y  $T_i$

$$E(Y_{1i} | T_i = 1, X_i) = E(Y_{1i} | T_i = 0, X_i)$$

$$E(Y_{1i} | T_i = 1, X_i) - E(Y_{0i} | T_i = 0, X_i) = \text{EFECTO CAUSAL PROMEDIO}$$

- Si no controlamos por  $X_i$ , se va al término de error
- Clave: determinar cuáles  $X_i$  incorporar

## Selección en observables

Tomemos la ecuación de la regresión

$$Y_i = \beta_0 + \beta_1 T_i + u_i$$

- Comparación de observables no ofrece buena medida pues tiene un sesgo de selección asociado

recordemos que los valores esperados de  $Y_i$  dependerán no solo del tratamiento si no que de más cosas

## Omisión de variable relevante

Poco frecuente que tengamos todas las variables que determinan un  $X_i$  tal que controlamos todos los factores contenido en los residuos.

- Si  $\text{corr}(u, X_i) \neq 0$  sesgo de selección
- Si  $E(u | T_i) \neq 0$

Imaginemos el modelo real

$$Y = X_1 \beta_1 + X_2 \beta_2 + u$$

Tenemos dos set de variables independientes ( $X_1$  y  $X_2$  vectores de variables)

El modelo que estimamos (excluye  $X_2$ )

$$Y = X_1 \tilde{\beta}_1 + u$$

$\tilde{\beta}_1$

$(X_1' X_1)^{-1} X_1' X_2 \beta_2$

SESGO DE VARIABLE OMITIDA

regresión de las variables en  $X_2$  parcializadas y luego presas en conjunto en matriz de todas las variables  $X_1$

Efecto de las  $X_2$  sobre  $Y$  (tenemos una intuición sobre  $\beta_2$ )

Esto solo será cero si  $X_1$  no está linealmente relacionado a las variables en  $X_2$  ( $X_1$  y  $X_2$  ortogonales)

$\beta_2 = 0$  (no significativo)

(A) RELACIÓN CON  $X_i$  correlación  $X_1$  y  $X_2$  omitida

(B) RELACIÓN CON  $Y_i$

Si alguna variable omitida cumple con [A] y [B]  $\rightarrow \tilde{\beta}_1$  estará sesgada

## DIRECCIONES DE SESGO

- Si  $\tilde{\beta}_1 > \beta_1$  (pensando  $\beta_2 +$ )  $\Rightarrow$  sesgo positivo correlación positiva } Modelo inicial sobre-estimo el efecto de  $X_1$
- Si  $\tilde{\beta}_1 < \beta_1$  (pensando  $\beta_2 +$ )  $\Rightarrow$  sesgo negativo correlación negativa } Modelo inicial subestimo el efecto real de  $X_1$
- Si  $\tilde{\beta}_1 < \beta_1$  (pensando  $\beta_2 -$ )  $\Rightarrow$  sesgo negativo correlación negativa } Modelo inicial subestimo el efecto real de  $X_1$
- Si  $\tilde{\beta}_1 > \beta_1$  (pensando  $\beta_2 -$ )  $\Rightarrow$  sesgo positivo correlación positiva } Modelo inicial sobre-estimo el efecto de  $X_1$

$\beta_{1,2}$  Pendientes de  $X_2$  sobre todas las variables explicativas

$\sim \text{corr}(X_1, X_2) \rightarrow$  Si están correlacionadas  $\oplus$  sesgo positivo  $\sim \beta_2 \oplus$  + intuición sobre signo  $\beta_2$   
 $\ominus$  sesgo negativo  $\sim \beta_2 \ominus$

Corr entre cualquier predictor en  $X_1$  y predictor omitido en  $X_2$ , será el vector de  $\beta_1$

Permite entender mejor al sesgo de selección

Si grupos difieren en características no observables,

MCO entregará estimación sesgada del impacto causal  $\rightarrow$  se ira al error

## Inclusión de variables irrelevantes

$\rightarrow$  último dígito del rui

Sobre-especificado cuando se incluyen variables que no forman parte del modelo poblacional

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_K X_K + \beta_{K+1} X_{K+1} + u$$

y esperamos que  $\beta_{K+1} = 0$

$\rightarrow E(\hat{\beta} - \beta) = 0$  dado que  $\neq$  INSESGADO

$\rightarrow$  Si  $X_{K+1}$  relacionado con  $X_1, \dots, X_K$  MULTICOLINEALIDAD  
 $\downarrow$   
varianzas de  $\hat{\beta}$

Incluso en un experimento aleatorio es mejor incluir variables control pese a que  $\neq$  duda  
potencialmente generan  $\neq$  mejoran la precisión al reducir varianza del error de  $\beta_1$

$$SE(\hat{\beta}_1) = \sqrt{\text{Var}(\hat{\beta}_1 | X)} = \frac{\sigma^2}{\sqrt{SCT_1 \cdot (1 - R_1^2)}}$$

Siempre va a crecer al incorporar variables toda vez que estén esas

- también más mediciones

cada vez que estén esas  
variables relacionadas