

## Primer ejercicio

Por motivo de la grave situación económica en el país luego de la pandemia, el Gobierno lanzó un programa social destinado a mejorar los ingresos de hogares pequeños cultivadores de café. En particular, la ayuda consistió en una capacitación acerca de la elección y uso de fertilizantes, ofrecida al jefe de hogar, junto con la entrega de una cantidad fija de dinero. Sin embargo, los recursos para desarrollar esta política eran limitados, por lo que sólo se podía brindar el programa a algunos de los hogares registrados en las bases de datos del Gobierno (los hogares no se tenían que inscribir para ser candidatos a recibir el programa). Por lo tanto, buscando evitar favorecimientos indebidos, se eligió aleatoriamente un subconjunto de hogares que recibirían el programa. Todos las familias seleccionadas para recibir la ayuda participaron activamente.

En busca de evaluar la efectividad de la medida, el Gobierno recolectó información acerca de la cosecha posterior al programa para todos los  $n$  hogares que fueron elegibles para el programa (esto es, tanto para los que recibieron el programa como para aquellos que no). Entre la información recolectada, se cuenta con la cantidad de kilogramos de café producido por hectárea cultivada por cada hogar  $C_i$ , así como una dummy  $D_i$  que toma el valor de uno si el hogar participó en el programa y cero de lo contrario.

Confundiendo en las habilidades de los estudiantes de econometría avanzada, a ustedes los contrata el Departamento Nacional de Planeación para llevar a cabo la evaluación de impacto respectiva.

Su jefe les dice que una manera de modelar el problema es a través de un modelo de regresión lineal dado por

$$C_i = \alpha + \tau D_i + \epsilon_i \quad (1)$$

donde  $\epsilon_i$  es un componente aleatorio idiosincrásico de media cero ( $\mathbb{E}[\epsilon_i] = 0$ ) y con varianza  $\sigma^2$  ( $\mathbb{E}[\epsilon_i^2] = \sigma^2$ ).

- a) Sean  $C_i(1)$  y  $C_i(0)$  los resultados potenciales de haber participado o no en el programa respectivamente. Similarmente, sean  $\epsilon_i(1)$  y  $\epsilon_i(0)$  los resultados potenciales análogos del error idiosincrático. Suponga además que  $\mathbb{E}[\epsilon_i(1)] = \mathbb{E}[\epsilon_i(0)] = 0$ .

- i) ¿Cuáles son las formas funcionales de  $C_i(1)$  y  $C_i(0)$  inducidas por el modelo (1)?

- Solución:

$$C_i(1) = \alpha + \tau + \epsilon_i(1)$$

$$C_i(0) = \alpha + \epsilon_i(0)$$



- ii) ¿Cuál parámetro del modelo captura el ATT del programa? Justifique matemáticamente.

- Solución:

El ATT se define como:  $\tau_{ATT} = (\bar{C}/D = 1) - (\bar{C}/D = 0)$ .

Sabemos que  $E[\epsilon_i/D_i] = 0$ , por lo que en el modelo de regresión  $\hat{\tau}$  es insesgado y consistente.

Entonces, tomando expectativa condicional del ATT definido y teniendo en cuenta el supuesto de **exogeneidad** se tiene que:

$$E[C_i(1)/D_i = 1] = E(\alpha + \tau + \epsilon_i/D_i = 1)$$

$$E[C_i(1)/D_i = 1] = \alpha + \tau + \cancel{E[\epsilon_i/D_i = 1]}$$

$$E[C_i(1)/D_i = 1] = \alpha + \tau$$

$$E[C_i(0)/D_i = 0] = E[\alpha + \epsilon_i/D_i = 0]$$

$$E[C_i(0)/D_i = 0] = \alpha + E[\epsilon_i/D_i = 0]$$

$$E[C_i(0)/D_i = 0] = \alpha$$

Restando las dos expresiones de arriba para tener al parámetro del  $\tau_{ATT}$ :

$$\tau_{ATT} = E[C_i(1)/D_i = 1] - E[C_i(0)/D_i = 0]$$

$$\tau_{ATT} = \alpha + \tau - \alpha$$

$$\tau_{ATT} = \tau$$

Por lo que  $\tau$  es el parámetro que captura del ATT en el modelo.

- iii) Imagine un escenario donde el programa, en lugar de ser asignado aleatoriamente, se entregaba a las familias que viven en lugares con climas menos favorables para el cultivo. ¿En este escenario el parámetro del inciso ii) sigue capturando el ATT?

**Pista:** Recuerde que

$$C_i = D_i C_i(1) + (1 - D_i) C_i(0); \quad \epsilon_i = D_i \epsilon_i(1) + (1 - D_i) \epsilon_i(0)$$

- Solución:

Dado que el tratamiento ya no es aleatorio, ahora va a haber un sesgo de selección, por lo que el parámetro  $\hat{\tau}$  ya no es estimador consistente de  $\tau_{ATT}$ :

$$\tau_{ATT} = E[C_i(1)/D_i] - E[C_i(0)/D_i]$$

Ya que no hay autorización:  $E[C_i(0)/D_i = 0] \neq E[C_i(0)/D_i = 1]$  entonces:

$$\tau_{ATT} = E[C_i(1) - C_i(0)/D_i = 1] + [E[C_i(0)/D_i = 1] - E[C_i(0)/D_i = 0]]$$

Con lo que se llega a que:

$$\tau_{ATT} = \hat{\tau} + E[C_i(0)/D_i = 1] - E[C_i(0)/D_i = 0]$$

Por lo que vemos que  $\hat{\tau}$  no captura el  $\tau_{ATT}$ .

- b) Bajo el cumplimiento de los supuestos del modelo clásico lineal, demuestre que el estimador de MCO  $\hat{\tau}$  es un estimador consistente del ATT y derive explícitamente su distribución asintótica (la varianza asintótica debe depender únicamente de  $n$ ,  $\sigma^2$  y  $p = \mathbb{P}(D_i = 1)$ ).

**Pista:** Recuerde el siguiente teorema visto en clase:

Una secuencia de vectores aleatorios  $\{x_N : N = 1, 2, \dots\}$  de  $K \times 1$  converge en distribución al vector aleatorio  $x$ , si para cualquier vector no aleatorio  $c$  de  $K \times 1$ ,

$$c^T x_N \xrightarrow{d} c^T x$$

- Solución:

Tenemos que por definición:

$$ATT = E[C_i(1)/D_i] - E[C_i(0)/D_i]$$

$$ATT = E[C_i(1) - C_i(0)/D_i = 1] + E[C_i(0)/D_i = 1] - E[C_i(0)/D_i = 0]$$

Y bajo supuestos MCL se tiene que:

$$ATT = \tau$$

Por lo que se que  $\tau = ATT$ , entonces  $\tau_{MCO}$  es consistente de  $\tau$  ? :

Para probarlo, primero sabemos que por definición y bajo supuestos MCL,  $\tau_{MCO}$  es:

$$\tau_{MCO} = \bar{C}_i(1)/D_i = 1 - \bar{C}_i(0)/D_i = 0$$

Aplicando  $plim$  a la expresión anterior se tiene:

$$plim(\tau_{MCO}) = plim(\bar{C}_i(1)/D_i = 1) - plim(\bar{C}_i(0)/D_i = 0)$$

Por WLLN tenemos que:

$$plim(\tau_{MCO}) = E[C_i(1)/D_i = 1] - E[C_i(0)/D_i = 0]$$

Y como ya se mostró, esta expresión bajo supuestos de MCL resulta en:

$$plim(\tau_{MCO}) = \tau$$

Por lo que se concluye que el  $\tau_{MCO}$  es estimador consistente del ATT.

Ahora para derivar la distribución asintótica: Sabemos que el modelo es  $C_i = \alpha + \tau D_i + \epsilon_i$  y sabemos que:

$$\beta \xrightarrow{a} N(\beta, \frac{1}{n} \sigma^2 Q_{xx}^{-1})$$

Donde  $\beta$  es el vector de estimadores del modelo de regresión  $(\alpha, \tau)$ . Entonces, sabiendo la distribución de  $\beta$  puedo saber la distribución de  $\tau$  usando el teorema y usando un vector  $C$  tal que la primera entrada sea cero y la segunda 1, de manera que:

$$\tau_{MCO} \xrightarrow{a} N(\tau, \frac{1}{n} \sigma^2 Q_{xx}^{-1})$$

Ahora bien, el término  $\sigma^2 Q_{xx}^{-1}$  no es observable, pero sé que  $\hat{\sigma}^2$  es estimador consistente de  $\sigma$  y  $(\frac{D_i' D_i}{n})^{-1}$  es estimador consistente de  $Q_{xx}^{-1}$ , y también sé que para la varianza de  $\tau$  solo necesito la entrada (2,2) de la matriz  $(\frac{D_i' D_i}{n})^{-1}$ .

La entrada (2,2) de  $(\frac{D_i' D_i}{n})^{-1}$  es :

$$\frac{\sum D_i' D_i}{n}$$

Dado que el vector de  $D_i$  es de 1s y 0s, por lo que  $\sum D_i' D_i$ , da como resultado simplemente la cantidad de 1s en la muestra, es decir la cantidad de tratados. Ahora eso dividido  $n$ , me da la probabilidad de ser tratado:  $p = \mathbb{P}(D_i = 1)$ .

Entonces, puedo reescribir la distribución de  $\tau$  así:

$$\tau_{MCO} \xrightarrow{a} N(\tau, \frac{1}{n} \hat{\sigma}^2 p)$$

Contentos por los reveladores hallazgos de los incisos a) y b), ustedes van por un tinto a la cafetería. En ella, encuentran dos colegas debatiendo acerca de cómo estimar correctamente el ATT. Uno de ellos argumenta que, debido a que hay aleatorización, el estimador debe ser una diferencia “ingenua” de medias de la variable dependiente entre los hogares participantes y aquellos que, aunque elegibles, no accedieron al programa. Su otro compañero lo contradice, pues afirma que la manera correcta de hacerlo es estimando  $\tau$  del modelo (1) por MCO dadas las conclusiones del inciso a) y b). ¿Quién tiene la razón?

Queriendo resolver esta encrucijada, a ustedes se les ocurre una idea: si logran probar que ambos estimadores son numéricamente equivalentes, entonces no debería importar cuál procedimiento utilicen.

- c) Prueben que la diferencia ingenua de medias es numéricamente equivalente al estimador de MCO del parámetro  $\tau$  del modelo (1).

**Ayuda:** Considere el modelo vectorial

$$C = X\beta + \epsilon$$

donde

$$\mathbf{X} = \begin{pmatrix} 1 & D_1 \\ 1 & D_2 \\ \vdots & \vdots \\ 1 & D_n \end{pmatrix}; \quad \mathbf{C} = \begin{pmatrix} C_1 \\ C_2 \\ \vdots \\ C_n \end{pmatrix}; \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}; \quad \beta = \begin{pmatrix} \alpha \\ \tau \end{pmatrix}$$

Puede usar sin probar que

$$(\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{\left(\sum_{i=1}^n D_i\right) \left(\sum_{i=1}^n (1 - D_i)\right)} \begin{pmatrix} \sum_{i=1}^n D_i & -\sum_{i=1}^n D_i \\ -\sum_{i=1}^n D_i & n \end{pmatrix}$$

y que

$$\mathbf{X}'\mathbf{C} = \begin{pmatrix} \sum_{i=1}^n C_i \\ \sum_{i=1}^n D_i C_i \end{pmatrix}$$

- Solución:  
Sé que:

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

Por lo que operando las matrices y sacando la entrada (2,1) de la matriz resultante tengo que:

$$\hat{\tau}_{MCO} = \frac{n \sum D_i C_i - \sum D_i \sum C_i}{\sum D_i (\sum (1 - D_i))}$$

Ahora bien, sé que cuando el modelo de regresión tiene una matriz X de tamaño (N x 2) y tiene una columna de 1s y una columna de Dummies, el estimador que acompañar la Dummy es la diferencia de medias simple:

$$C_i = \alpha + \tau D_i + \epsilon_i$$

$$\hat{\tau}_{MCO} = \bar{C}_1 - \bar{C}_0$$

Entonces, noto que:

$$\bar{C}_1/D_i = 1 = \frac{\sum C_i D_i}{\sum D_i} = \frac{\sum C_i}{\sum 1} = \frac{\sum C_1}{n}$$

$$\bar{C}_0/D_i = 0 = \frac{\sum C_i (1 - D_i)}{\sum (1 - D_i)} = \frac{\sum C_i}{\sum 1 - 0} = \frac{\sum C_0}{n}$$

Ahora, defino la diferencia ingenua de media como:

$$\hat{\tau}_{ing} = \bar{C}_1/D_i = 1 - \bar{C}_0/D_i = 0$$

Ahora reemplazando, tengo que:

$$\hat{\tau}_{ing} = \frac{\sum C_i D_i}{\sum D_i} - \frac{\sum C_i (1 - D_i)}{\sum (1 - D_i)}$$

$$= \frac{\sum (1 - D_i) \sum D_i C_i - \sum D_i \sum C_i (1 - D_i)}{\sum D_i (\sum (1 - D_i))}$$

Operando solo el numerador, tengo que puedo cancelar las  $\sum D_i$  cuando  $D_i = 0$  y se que  $\sum 1 = n$ :

$$(n - \cancel{\sum D_i}) \sum C_i D_i - \sum D_i (\sum C_i (1 - \cancel{D_i}))$$

Con lo que llego a la expresión:

$$\hat{\tau}_{ing} = \frac{n \sum D_i C_i - \sum D_i \sum C_i}{\sum D_i (\sum (1 - D_i))} = \hat{\tau}_{MCO}$$

Se demuestra entonces que  $\hat{\tau}_{ing} = \hat{\tau}_{MCO}$

Ahora una compañera le comenta que, según lo que ha leído, ella cree que el éxito del programa depende en gran medida de la calidad de la tierra donde se siembra. De manera que, pueden existir posibles *efectos heterogéneos* dependiendo de las dotaciones de este factor en cada hogar. Asimismo, le comenta que afortunadamente se cuenta con una variable  $Z_i$  en la base de datos que captura esta información. Por lo tanto, ella propone que un modelo más completo estaría dado por

$$C_i = \beta_0 + \beta_1 D_i + \beta_2 \tilde{Z}_i + \beta_3 D_i \tilde{Z}_i + \epsilon_i \quad (2)$$

donde  $\tilde{Z}_i = (Z_i - \bar{Z})$  y  $\bar{Z}$  es el promedio muestral. Finalmente, suponga que  $\mathbb{E}[\epsilon_i(1)|Z_i] = 0$  y  $\mathbb{E}[\epsilon_i(0)|Z_i] = 0$ , esto es, que la calidad de la tierra es exógena.

d) Partiendo del modelo (2):

i) Calcule el efecto esperado en los tratados como función de  $Z$ :

$$ATT(Z) = \mathbb{E}[C_i(1) - C_i(0)|D_i = 1, Z]$$

¿Cuál es la interpretación de  $\beta_1$  y  $\beta_3$ ?

• Solución:

$$ATT(Z) = E[C_i(1) - C_i(0)|D_i = 1, Z]$$

Aplicando esperanza condicional:

$$E[C_i(1)|D_i = 1, Z] = \beta_0 + \beta_1 + \beta_2 E[\hat{Z}|D_i = 1] + \beta_3 E[\hat{Z}|D_i = 1] + \cancel{E[\epsilon_i|D_i = 1, Z]}$$

(Se cancela el valor esperado condicional del error por supuesto de exogeneidad), y por la aleatorización se que  $E[C_i(0)|D_i = 1, Z] = E[C_i(0)|D_i = 0, Z]$

$$E[C_i(0)|D_i = 1, Z] = \beta_0 + \beta_2 E[\hat{Z}|D_i = 0] + \cancel{E[\epsilon_i|D_i = 0, Z]}$$

Uniendo los términos:

$$ATT(Z) = \cancel{\beta_0} + \beta_1 + \cancel{\beta_2 E[\hat{Z}|D_i = 1]} + \beta_3 E[\hat{Z}|D_i = 1] - \cancel{\beta_0} - \cancel{\beta_2 E[\hat{Z}|D_i = 0]}$$

Llego a que:

$$ATT(Z) = \beta_1 + \beta_3 E[\hat{Z}|D_i = 1]$$

Usando LEI:

$$ATT(Z) = \beta_1 + \beta_3 E[\hat{Z}]$$

En este escenario  $\beta_1$  se interpreta como el cambio esperado en  $C_i$  cuando cambio  $D_i$ , dejando todo lo demás constante. En última el  $\beta_1$  es el **efecto del tratamiento**. Por otro lado,  $\beta_3$  se puede entender como las desviaciones del efecto del tratamiento.

ii) Halle el parámetro que ahora captura el ATT del programa.

• Solución: Sabemos que :

$$ATT = \beta_1 + \beta_3 \hat{Z}$$

$$ATT(Z) = \beta_1 + \beta_3 E[\hat{Z}]$$

Sé que por definición:  $E[\hat{Z}] = E[Z_i - \bar{Z}] = 0$ , por lo que llego a:

$$ATT = \beta_1$$

Entonces, el  $\beta_1$  captura el ATT en este modelo.

iii) ¿Qué pasaría si  $Z_i$  no fuese exógena?

• Solución:

Si  $Z_i$  no es exógena entonces  $E[\epsilon_i|Z_i] \neq 0$  y por ende el estimador MCO no sería consistente ni insesgado.

Finalmente, su jefe les tiene un último reto: probar si para al menos el 95 % de los hogares participantes el efecto de participar fue positivo.

e) Suponga que  $Z_i \sim N(\mu, 1)$ . Usando la forma funcional del modelo (2):

i) Encuentre la distribución de  $ATT(Z)$

- Solución:

$$ATT(Z) = \beta_1 + \beta_3 E[\hat{Z}]$$

Dado que  $Z_i \sim N(\mu, 1)$  y que  $\hat{Z}_i = (Z_i - \bar{Z})$ , sé que:

$$\hat{Z}_i \sim N(0, 1)$$

Sabiendo la distribución de  $\hat{Z}_i$ , puedo saber la distribución de  $ATT(Z)$ :

$$ATT(Z) \sim N(\hat{\beta}_1, \hat{\beta}_3^2)$$

- ii) Encuentre el percentil 5 ( $p_{0.05}$ ) de la distribución de  $ATT(Z)$  en términos de los parámetros del modelo.

- Solución:

$$P(ATT(Z) \leq p_{0.05}) = 0.05$$

Estandarizando:

$$P\left(\frac{ATT(Z) - \beta_1}{\beta_3} \leq \frac{p_{0.05} - \beta_1}{\beta_3}\right) = 0.05$$

Sé que:

$$\frac{ATT(Z) - \beta_1}{\beta_3} \sim N(0, 1)$$

Cómo es una normal estándar, el punto crítico de la distribución en el percentil 5 es 1,64:

$$P(1,64 \leq \frac{p_{0.05} - \beta_1}{\beta_3}) = 0.05$$

Por último, despejando tengo que:

$$\hat{p}_{0.05} = 1,64\hat{\beta}_3 + \hat{\beta}_1$$

- iii) Suponga que usted estimó por MCO que

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{pmatrix} = \begin{pmatrix} 0.6 \\ 2 \\ 0.3 \\ -0.03 \end{pmatrix}; \quad S = \begin{pmatrix} 0.6 & -0.26 & 0.33 & 0 \\ -0.26 & 0.7 & 0.5 & 0 \\ 0.33 & 0.5 & 0.35 & -0.25 \\ 0 & 0 & -0.25 & 0.08 \end{pmatrix}$$

donde  $S$  es la matriz de varianza-covarianza estimada de los parámetros.

Construya un estimador consistente de  $p_{0.05}$  y diseñe un test que le permita probar:

$$\begin{cases} H_0 : p_{0.05} \leq 0 \\ H_a : p_{0.05} > 0 \end{cases}$$

Use los datos disponibles para ejecutar su test. Concluya para un nivel de significancia  $\alpha = 0.05$ .

- Solución:

$$\hat{p}_{0.05} = 1,64\hat{\beta}_3 + \hat{\beta}_1$$

Puedo saber la distribución de  $\hat{p}_{0.05}$ :

$$\hat{p}_{0.05} \sim N(1,64\hat{\beta}_3 + \hat{\beta}_1, (1,64)^2 var(\hat{\beta}_3) + var(\hat{\beta}_1) + 2(1,64)cov(\hat{\beta}_1, \hat{\beta}_3))$$

De ahí puedo tener un estadístico teórico  $Z_c$ :

$$Z_c = \frac{\hat{p}_{0.05} - p_{0.05}}{\sqrt{var(\hat{p}_{0.05})}} = \frac{1,64\hat{\beta}_3 + \hat{\beta}_1 - (1,64\beta_3 + \beta_1)}{\sqrt{(1,64)^2 var(\hat{\beta}_3) + var(\hat{\beta}_1) + 2(1,64)cov(\hat{\beta}_1, \hat{\beta}_3)}}$$

De lo anterior, se tiene que  $Z_c \sim N(0, 1)$  dado que está estandarizado. Con esto puedo construir el estadístico de prueba  $t_c$ :

$$t_c = \frac{\hat{p}_{0.05} - p_{0.05}}{\sqrt{var(\hat{p}_{0.05})}} = \frac{(1,64)(-0.03) + 2 - (1,64\beta_3 + \beta_1)}{\sqrt{(1,64)^2(0,08) + (0,7) + 2(1,64)0}}$$

Bajo  $H_0$  se tiene que:  $\beta_1 = \beta_3 = 0$ , por lo que la prueba sería:

$$t_c = \frac{(1,64)(-0.03) + 2}{\sqrt{(1,64)^2(0,08) + (0,7)}} = 2,039$$

A un nivel de significancia del 0.05, **se rechaza la hipótesis nula en favor de la alterna.**

- iv) Verdadero o falso: Si el  $p$ -valor del test es mayor a 0.05, entonces existe evidencia estadística de que el efecto de que el hogar en el percentil 5 de la distribución de efectos fue negativo o nulo. Justifique.



## Segundo ejercicio

Una curva de aprendizaje se puede pensar como el cambio en la productividad o eficiencia con la que se hacen las cosas en el transcurso del tiempo. Este concepto es utilizado en Economía, por ejemplo, para modelar el cambio de los costos reales de las firmas según el tiempo que llevan operando. En particular, se supone que las firmas van aprendiendo a producir de manera más eficiente, por lo que el costo real de producción debería, *ceteris paribus*, decrecer.

Supongan que la curva generalizada de aprendizaje de una firma hipotética está dada por

$$C_i = C_0 N_i^{\frac{\alpha}{\gamma}} Y_i^{\frac{1-\alpha}{\gamma}} \exp(u_i), \quad (1)$$

donde  $C_i$  corresponde a los costos reales unitarios que enfrenta una firma  $i$ ;  $Y_i$  es su nivel de producción;  $N_i$  es la producción acumulada a lo largo del tiempo;  $C_0$  corresponde a una medida de costos iniciales; y  $u_i$  es un término estocástico de media cero desconocido para el econometrista. El parámetro  $\alpha$  determina la dirección de la elasticidad del costo unitario con respecto a la producción acumulada. Este es el principal parámetro de interés. Finalmente, el parámetro  $\gamma \in \mathbb{R}^+$  caracteriza los retornos a escala de la función de aprendizaje: si  $\gamma = 1$  la curva tiene retornos constantes a escala, si  $\gamma < 1$  la curva tiene retornos decrecientes a escala, y si  $\gamma > 1$  la curva tiene retornos crecientes a escala.<sup>1</sup>

Usted, como investigador, está interesado en estimar la curva de aprendizaje que enfrentan las firmas de la industria manufacturera colombiana. Para ello, cuenta con información de  $C_i$ ,  $N_i$  y  $Y_i$  para  $N$  firmas de la industria.

- a) Propongan un modelo de regresión lineal que capture la forma funcional dada en (1). Muestre cómo los parámetros de su nuevo modelo dependen de los parámetros de la curva de aprendizaje. Mencionen y discutan los supuestos necesarios para que los estimadores por MCO de los parámetros del modelo sean insesgados y consistentes. ¿Son estos supuestos plausibles en el contexto del problema?

■ Solución:

La ecuación 1 representa la *curva de aprendizaje* de que enfrentan las firmas de la industria manufacturera colombiana. Se observa, que la relación entre los costos reales unitarios que enfrenta una firma  $i$  ( $C_i$ ), respecto a su nivel de producción ( $Y_i$ ) y la producción acumulada a lo largo del tiempo ( $N_i$ ) es no lineal. No obstante, por la forma de la curva de aprendizaje, si se aplica el logaritmo natural como transformación matemática a la curva de aprendizaje original, es posible representar la relación entre las variables de interés del modelo como un modelo de regresión lineal<sup>2</sup> que capture la forma funcional de 1.

Al aplicar el logaritmo natural a la curva de aprendizaje original se obtiene el siguiente *modelo de regresión lineal*:

$$\log(C_i) = \beta_0 + \beta_n \log(N_i) + \beta_y \log(Y_i) + u_i \quad (3)$$

Donde, claramente se observa que los parámetros del nuevo modelo dependen de los parámetros de la curva de aprendizaje:

- $\beta_0 = \log(C_0)$

<sup>1</sup>Si la tecnología de producción muestra rendimientos constantes a escala, los costos unitarios reales no deben variar con el nivel de producción. Por el contrario, si los rendimientos son crecientes a escala, los costos unitarios deberían disminuir a medida que aumenta el nivel de producción, y si son decrecientes, se esperaría lo contrario.

<sup>2</sup>Básicamente que se pueda representar la curva de aprendizaje como una relación matemática lineal en los parámetros del modelo.

- $\beta_n = \frac{\alpha}{\gamma}$
- $\beta_y = \left( \frac{1-\alpha}{\gamma} \right)$

Ahora bien, el *modelo clásico lineal (MCL)* tiene cuatro supuestos básicos:

1. El modelo es lineal en los parámetros:

$$y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{ik}\beta_k + \varepsilon_i$$

2. Exogeneidad de las variables independientes:

$$E[\varepsilon_i | x_{j1}, x_{j2}, \dots, x_{jk}] = 0 \quad \text{para } j = 1, \dots, N$$

El cual entre otras cosas es un supuestos clave para la identificación de parámetros en el MCL.

3.  $\mathbf{X}$  es una matriz estocástica y de rango completo:

$$\text{rank}(\mathbf{X}) = K$$

El cual también es un supuesto clave para la identificación de parámetros en el MCL.

4. Homocedasticidad y no autocorrelación:

$$\begin{aligned} \text{Var}[\varepsilon_i | \mathbf{X}] &= \sigma^2 \quad \text{para } i = 1, \dots, N \\ \text{Cov}[\varepsilon_i, \varepsilon_j | \mathbf{X}] &= 0 \quad \text{para } i \neq j \end{aligned}$$

De los 4 supuestos del MCL mencionados anteriormente, solo los 3 primeros son necesarios para que los estimadores por MCO de los parámetros del modelo sean *insesgados* y *consistentes*. El *supuesto de linealidad en los parámetros* es necesario para poder escribir el modelo en términos de operadores lineales representados por las matrices que hacen parte de la formulación del MCL y del estimador OLS. El *supuesto de exogeneidad* es necesario para garantizar que  $E[\varepsilon | \mathbf{X}] = 0$ , lo que garantiza que  $E[\hat{\beta}_{MCO}] = \beta$ . Finalmente, el *supuesto de rango completo* es necesario para que  $\mathbf{X}'\mathbf{X}$  sea invertible.

Nota: Importante mencionar que en ningún momento se usa el supuestos de homocedasticidad ni correlación serial para garantizar insesgamiento o consistencia en el estimador MCO.

Los 3 supuestos necesarios para que el estimador MCO sea insesgado y consistente son supuestos plausibles en el contexto del problema:

- *El supuestos de linealidad en los parámetros*: es plausible en la medida de que una curva de aprendizaje tipo *Cobb-Douglass* como la planteada en la ecuación 1 captura los elementos más importantes en nivel de producción ( $Y_i$ ) y producción acumulada a lo largo del tiempo ( $N_i$ ) para explicar el costo real unitario ( $C_i$ ). De igual forma, también tiene sentido introducir el error de forma  $\exp(u_i)$  dado que es un reescalamiento útil de un término estocástico  $u_i$  de media cero y desconocido para el investigador. Como ya se explico antes, mediante una transformación matemática logarítmica es posible escribir la curva de costo original 1 como un modelo de regresión lineal en los parámetros como el que muestra el modelo 3.
- *El supuestos de exogeneidad*: también es posible en la medida que las variables más relevante a la hora de explicar el costo real unitario ( $C_i$ ) posiblemente son el nivel de producción ( $Y_i$ ) y la producción acumulada a lo largo del tiempo ( $N_i$ ). Si la información relevante para explicar  $C_i$  se encuentran principalmente en  $Y_i$  y en  $N_i$ , es muy factible que cualquier componente idiosincrático no observable que se encuentre almacenado en  $u_i$  sea exógeno o no tenga relaciones lineal o no lineal con  $Y_i$  y  $N_i$ , es decir se satisface que  $E[\varepsilon_i | Y_i, N_i] = 0$ .
- *El supuestos de rango completo de  $\mathbf{X}$* : también es un supuesto muy plausible en la medida que no es factible realmente que haya un problema de multicolinealidad perfecta en el modelo de regresión lineal 3, dado que no es factible que haya una combinación lineal perfecta entre los regresores  $Y_i$  y  $N_i$ , lo que implica que las columnas de la matriz  $\mathbf{X}$  son linealmente independientes.

No obstante, se resalta que el supuesto de *homocedasticidad y no autocorrelación* es muy poco plausible que se cumpla dado que se espera que hayan elementos idiosincráticos no observables que hagan parte del componente  $u_i$  de la curva de aprendizaje de cada firma que afecten la dispersión o varianza de dicho término  $u_i$ . Por tanto, es muy factible que haya heterocedasticidad en la muestra de corte transversal de las firmas de la industria manufacturera colombiana y el supuesto de homocedasticidad y no autocorrelación no se satisfaga. No obstante, como dicho supuestos no se necesita para que el estimador de MCO sea *insesgado* y *consistente*, se tiene que si se estima la ecuación 3 por MCO los parámetros  $\beta_0 = \log(C_0)$ ,  $\beta_n = \frac{\alpha}{\gamma}$  y  $\beta_y = \left( \frac{1-\alpha}{\gamma} \right)$  van a ser insesgados y consistentes.



- b) Supongan que se cumplen los supuestos que ustedes discutieron en el inciso anterior. Propongan un estimador consistente de  $\alpha$  y  $\gamma$ .

■ Solución:

Lo primero que hay que notar es que el modelo que se va a estimar es el *modelo de regresión lineal* planteado en la ecuación 3. Los parámetros de interés que van a ser estimados en este modelo son  $\beta_n$  y  $\beta_y$ . De la discusión anterior, se sabe que si estima la ecuación 3 por MCO los estimadores que se obtengan para  $\beta_n$  y  $\beta_y$  son insesgados y consistentes, dado que, como se menciona en el literal anterior, es muy plausible que los 3 supuestos de: 1) linealidad en los parámetros, 2) exogeneidad  $E[u_i|Y_i, N_i] = 0$  y 3) rango completo de  $\mathbf{X}$  se satisfagan para el modelo de regresión lineal planteado en 3 y las características del problema que se está investigando. Entonces, al estimar 3 por MCO, se tiene que por consistencia de dicho estimador:

$$\text{plim}(\hat{\beta}_n) = \beta_n$$

$$\text{plim}(\hat{\beta}_y) = \beta_y$$

Ahora bien, se busca *estimadores consistentes* para los parámetros  $\alpha$  y  $\gamma$  que hacen parte de la curva de aprendizaje original de la ecuación 1.

Para ello, va a ser necesario emplear el *teorema de mapeo continuo o teorema de Slutsky* que dice que:

Sea  $\mathbf{g} : R^K \rightarrow R^J$  una función continua. Sea  $\{\mathbf{x}_N : N = 1, 2, \dots\}$  una secuencia de vectores aleatorios de  $K \times 1$  tal que  $\mathbf{x}_N \xrightarrow{p} \mathbf{c}$ , entonces  $\mathbf{g}(\mathbf{x}_N) \xrightarrow{p} \mathbf{g}(\mathbf{c})$

Como el modelo que se va a estimar es el modelo de regresión lineal 3, se van a obtener los estimadores MCO de  $\hat{\beta}_n$  y  $\hat{\beta}_y$ .

Se sabe que:

$$\beta_n = \frac{\alpha}{\gamma}$$

$$\beta_y = \left( \frac{1 - \alpha}{\gamma} \right)$$

Entonces, para recuperar los parámetros  $\alpha$  y  $\gamma$ , es necesario despejarlos de las dos ecuaciones anteriores.

EL parámetro  $\gamma$  se puede escribir como  $\gamma = \frac{\alpha}{\beta_n}$

Lo que hace que,

$$\beta_y = \frac{1 - \alpha}{\gamma}$$

$$\beta_y = \frac{1}{\gamma} - \frac{\alpha}{\gamma}$$

$$\beta_y = \frac{1}{\left(\frac{\alpha}{\beta_n}\right)} - \beta_n$$

$$\beta_y = \frac{\beta_n}{\alpha} - \beta_n$$

$$\frac{\beta_n}{\alpha} = \beta_n + \beta_y$$

$$\alpha = \frac{\beta_n}{\beta_n + \beta_y}$$

y

$$\gamma = \frac{\alpha}{\beta_n}$$

$$\gamma = \left( \frac{1}{\beta_n} \right) \left( \frac{\beta_n}{\beta_n + \beta_y} \right)$$

$$\gamma = \frac{1}{\beta_n + \beta_y}$$

Se proponen como estimadores consistentes de  $\alpha$  y  $\gamma$ :

$$\hat{\alpha} = \frac{\hat{\beta}_n}{\hat{\beta}_n + \hat{\beta}_y}$$

$$\hat{\gamma} = \frac{1}{\hat{\beta}_n + \hat{\beta}_y}$$

Si se define la función  $\mathbf{g}$  como:

$$\mathbf{g} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$\begin{pmatrix} \beta_n \\ \beta_y \end{pmatrix} \rightarrow \mathbf{g} \begin{pmatrix} \beta_n \\ \beta_y \end{pmatrix} = \begin{pmatrix} \alpha \\ \gamma \end{pmatrix} = \begin{pmatrix} \frac{\beta_n}{\beta_n + \beta_y} \\ \frac{1}{\beta_n + \beta_y} \end{pmatrix}$$

Ahora bien, si se aplica el *teorema de Slutsky* se tiene que:

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\gamma} \end{pmatrix} = \mathbf{g} \begin{pmatrix} \hat{\beta}_n \\ \hat{\beta}_y \end{pmatrix} \xrightarrow{p} \mathbf{g} \begin{pmatrix} \beta_n \\ \beta_y \end{pmatrix} = \begin{pmatrix} \alpha \\ \gamma \end{pmatrix} \quad (4)$$

Dado que, por consistencia del estimador de MCO al estimar la regresión lineal 3, se sabe que:

$$\begin{pmatrix} \hat{\beta}_n \\ \hat{\beta}_y \end{pmatrix} \xrightarrow{p} \begin{pmatrix} \beta_n \\ \beta_y \end{pmatrix}$$

Lo cual hace que  $\hat{\alpha}$  y  $\hat{\gamma}$  sean estimadores consistentes de  $\alpha$  y  $\gamma$ .

- c) Propongan estimadores de las desviaciones estándar de  $\hat{\alpha}$  y  $\hat{\gamma}$ . Para esto, supongan que la matriz de varianzas y covarianzas asintótica de  $\sqrt{N}(\hat{\beta} - \beta)$  y su respectivo estimador están dados por

$$V = Avar(\sqrt{N}(\hat{\beta} - \beta)) = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} \\ \sigma_{12} & \sigma_2^2 & \sigma_{23} \\ \sigma_{13} & \sigma_{23} & \sigma_3^2 \end{pmatrix} \quad \text{y} \quad \hat{V} = \widehat{Avar}(\sqrt{N}(\hat{\beta} - \beta)) = \begin{pmatrix} \hat{\sigma}_1^2 & \hat{\sigma}_{12} & \hat{\sigma}_{13} \\ \hat{\sigma}_{12} & \hat{\sigma}_2^2 & \hat{\sigma}_{23} \\ \hat{\sigma}_{13} & \hat{\sigma}_{23} & \hat{\sigma}_3^2 \end{pmatrix} \quad (5)$$

respectivamente. Asegúrense de que sus expresiones estén en términos de  $\hat{\beta}_1$ ,  $\hat{\beta}_2$  y los elementos que componen  $\hat{V}$ .

**Pista:** El método delta puede resultar útil.

■ Solución:

Podemos usar el *método delta* y el *teorema de Slutsky* para obtener la *matriz var-cov* y de ahí tener estimadores de las desviaciones estándar. En particular sabemos que una expresión de la forma:

$$J[\hat{V}]J^T$$

Donde  $J$  es la matriz de derivadas parciales<sup>3</sup> de los estimadores de  $\alpha$  y  $\gamma$ , da la matriz var-cov que necesitamos<sup>4</sup>.

Definimos  $J$ :

$$J = \begin{pmatrix} \frac{\partial \alpha}{\partial \beta_0} & \frac{\partial \alpha}{\partial \beta_1} & \frac{\partial \alpha}{\partial \beta_2} \\ \frac{\partial \gamma}{\partial \beta_0} & \frac{\partial \gamma}{\partial \beta_1} & \frac{\partial \gamma}{\partial \beta_2} \end{pmatrix} = \begin{pmatrix} 0 & \frac{\hat{\beta}_2}{(\hat{\beta}_1 + \hat{\beta}_2)^2} & \frac{-\hat{\beta}_1}{(\hat{\beta}_1 + \hat{\beta}_2)^2} \\ 0 & \frac{-1}{(\hat{\beta}_1 + \hat{\beta}_2)^2} & \frac{-1}{(\hat{\beta}_1 + \hat{\beta}_2)^2} \end{pmatrix}$$

Y ya sabemos que  $\hat{V}$  es:

$$\hat{V} = \begin{pmatrix} \hat{\sigma}_1^2 & \hat{\sigma}_{12} & \hat{\sigma}_{13} \\ \hat{\sigma}_{12} & \hat{\sigma}_2^2 & \hat{\sigma}_{23} \\ \hat{\sigma}_{13} & \hat{\sigma}_{23} & \hat{\sigma}_3^2 \end{pmatrix}$$

<sup>3</sup>También conocida en la literatura como matriz Jacobiana

<sup>4</sup>Por simplicidad de notación,  $\beta_n = \beta_1$  y  $\beta_y = \beta_2$

Entonces aplicando Slutsky tenemos:

$$\begin{pmatrix} 0 & \frac{\widehat{\beta}_2}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} & \frac{-\widehat{\beta}_1}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} \\ 0 & \frac{-1}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} & \frac{-1}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} \end{pmatrix} \begin{pmatrix} \widehat{\sigma}_1^2 & \widehat{\sigma}_{12} & \widehat{\sigma}_{13} \\ \widehat{\sigma}_{12} & \widehat{\sigma}_2^2 & \widehat{\sigma}_{23} \\ \widehat{\sigma}_{13} & \widehat{\sigma}_{23} & \widehat{\sigma}_3^2 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ \frac{\widehat{\beta}_2}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} & \frac{-1}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} \\ \frac{-\widehat{\beta}_1}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} & \frac{-1}{(\widehat{\beta}_1 + \widehat{\beta}_2)^2} \end{pmatrix}$$

Finalmente, operando las matrices se llega a:

$$\begin{pmatrix} \frac{\widehat{\beta}_2}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_2^2 \widehat{\beta}_2 - \widehat{\sigma}_{23} \widehat{\beta}_1}{\widehat{\phi}} \right] - \frac{\widehat{\beta}_1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_{23} \widehat{\beta}_2 - \widehat{\sigma}_3^2 \widehat{\beta}_1}{\widehat{\phi}} \right] & -\frac{1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_2^2 \widehat{\beta}_2 - \widehat{\sigma}_{23} \widehat{\beta}_1}{\widehat{\phi}} \right] - \frac{1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_{23} \widehat{\beta}_2 - \widehat{\sigma}_3^2 \widehat{\beta}_1}{\widehat{\phi}} \right] \\ \frac{\widehat{\beta}_2}{\widehat{\phi}} \left[ \frac{-\widehat{\sigma}_2^2 - \widehat{\sigma}_{23}}{\widehat{\phi}} \right] + \frac{\widehat{\beta}_1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_3^2 + \widehat{\sigma}_{23}}{\widehat{\phi}} \right] & \frac{1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_2^2 + \widehat{\sigma}_{23}}{\widehat{\phi}} \right] + \frac{1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_3^2 + \widehat{\sigma}_{23}}{\widehat{\phi}} \right] \end{pmatrix}$$

Donde  $\widehat{\phi} = (\widehat{\beta}_1 + \widehat{\beta}_2)^2$

Entonces, los estimadores de las desviaciones estándar  $\widehat{\alpha}$  y  $\widehat{\gamma}$  serán:

$$\widehat{DE}_{\widehat{\alpha}} = \sqrt{\frac{\widehat{\beta}_2}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_2^2 \widehat{\beta}_2 - \widehat{\sigma}_{23} \widehat{\beta}_1}{\widehat{\phi}} \right] - \frac{\widehat{\beta}_1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_{23} \widehat{\beta}_2 - \widehat{\sigma}_3^2 \widehat{\beta}_1}{\widehat{\phi}} \right]}$$

$$\widehat{DE}_{\widehat{\gamma}} = \sqrt{\frac{1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_2^2 + \widehat{\sigma}_{23}}{\widehat{\phi}} \right] + \frac{1}{\widehat{\phi}} \left[ \frac{\widehat{\sigma}_3^2 + \widehat{\sigma}_{23}}{\widehat{\phi}} \right]}$$

Usted cuenta con la base de datos “manufacturaCol.dta” para implementar la estimación de los parámetros propuesta en los incisos anteriores. La base de datos cuenta con información de las siguientes variables para cada una de las 11,500 firmas de la muestra:

- *costos*: Costos unitarios reales de producción que enfrenta la firma en el momento que se levantaron los datos medido en miles de millones pesos.
- *producto*: Producción de la firma en miles de millones pesos en el momento que se levantaron los datos.
- *producto\_acum*: Producción acumulada histórica de la firma en el momento en que se levantaron los datos medida en miles de millones de pesos .

Resuelvan los siguientes incisos a partir de los datos disponibles en la base de datos.

- d) A manera de estadísticas descriptivas, representen mediante un gráfico de dispersión y su respectiva línea de ajuste las siguientes relaciones: 1. La relación de  $\log(C_i)$  con  $\log(N_i)$  una vez se ha “removido” el efecto de  $\log(Y_i)$ . 2. La relación de  $\log(C_i)$  con  $\log(Y_i)$  una vez se ha “removido” el efecto de  $\log(N_i)$ .

**Ayuda:** para lograr esto, apliquen la intuición de partialling out o residualización de los modelos de regresión múltiples. Interprete los gráficos de dispersión.

- Solución:

Para la realización de las gráficas de dispersión y su respectiva línea de ajuste, se empleo el siguiente código:

```
1
2 // Taller 1
3 use "/Users/federicoduenas/Desktop/Econometri a_ Avz_/taller 1/manufacturaCol.dta"
4
5 foreach var of varlist producto_acum producto costos {
6     gen ln_var = ln('var'+1)
7 }
8
9
10 // relacion de log_c con log_n controlando por log_y
11 // usando Partialling out:
12
13 reg ln_producto_acum ln_producto
14
15 gen aux_1 = 9.256792 + 0.1120002*ln_producto
16
17 // para tener la parte de N que no es explicada por Y
18 gen e_1 = ln_producto_acum - aux_1
19
20 // ahora para tener los C_i no explicada por Y_i
21
```

```

22 reg ln_costos ln_producto
23
24 gen aux_2_1 = -.019997 + .0052894*ln_producto
25
26 // los residuales: parte de C_i no explicada por Y_i
27 gen e_1_2 = ln_costos - aux_2_1
28
29 // sactter de la relaci n C_i con N_i una vez se removi el efecto de Y_i
30 ***#
31
32 scatter e_1_2 e_1 || lfit e_1_2 e_1
33
34 /*
35
36
37 // ahora si calculo la relacion entre N y C, controlando por Y
38
39 reg ln_costos e_1
40
41 // Que, por teorema de Waugh-Frisch-Lovell, es equivalente al par metro de log_N de la reg
    lineal:
42 reg ln_costos ln_producto_acum ln_producto
43
44 // gr fica de la relaci n de C y N:
45 scatter ln_costos e_1 || lfit ln_costos e_1
46 */
47 // ahora la relacion de log_C y log_Y controlando por log_N
48
49 reg ln_producto ln_producto_acum
50
51 gen aux_2 = .3859203 + .4743226*ln_producto_acum
52
53 gen e_2 = ln_producto - aux_2
54
55
56 // ahora remuevo el efecto de N_i sobre C_i
57
58 reg ln_costos ln_producto_acum
59 gen aux_2_2 = -.013807 + .0020865*ln_producto_acum
60 gen e_2_2 = ln_costos - aux_2_2
61
62
63 // sactter de la relaci n de C y Y removiendo el efecto de N_i
64
65 scatter e_2_2 e_2 || lfit e_2_2 e_2
66
67 /*
68 reg ln_costos e_2
69
70 reg ln_costos ln_producto_acum ln_producto
71
72
73 scatter ln_costos e_2 || lfit ln_costos e_2

```

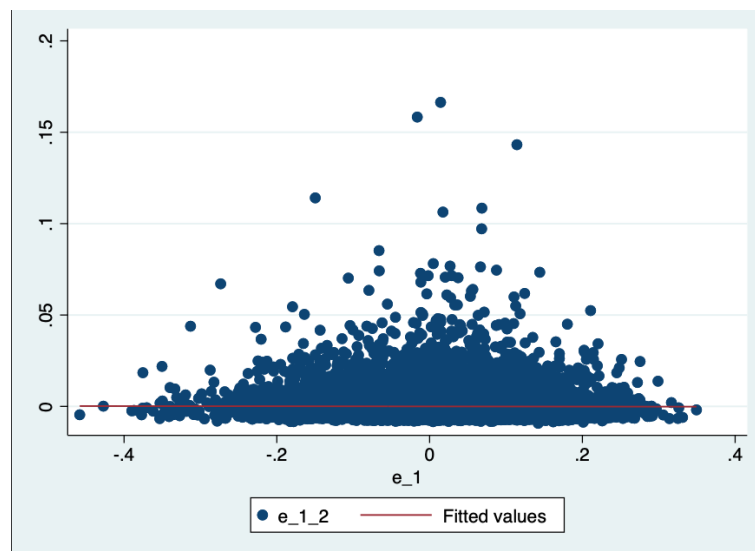


Figura 1: relación de  $\log(C_i)$  con  $\log(N_i)$  una vez se ha “removido” el efecto de  $\log(Y_i)$

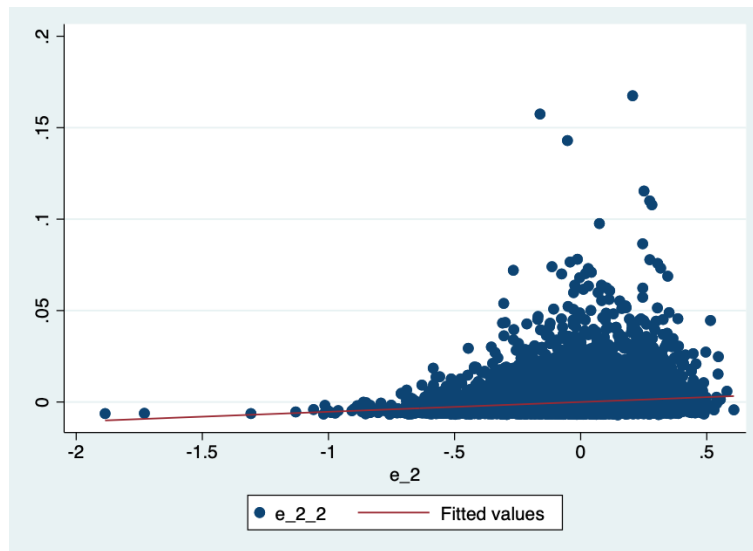


Figura 2: relación de  $\log(C_i)$  con  $\log(Y_i)$  una vez se ha “removido” el efecto de  $\log(N_i)$

De los gráficos de dispersión, se concluye que hay una baja asociación lineal o correlación entre  $\log(C_i)$  y  $\log(N_i)$  que se visualiza claramente por la alta dispersión de la gráfica de dispersión sin un patrón lineal claro y por lo relativamente plano que resulta ser la línea de ajuste asociada. Por otro lado, se observa que hay existe un poco más grado de asociación lineal o correlación entre  $\log(C_i)$  y  $\log(Y_i)$  (aunque también se podría decir que es relativamente bajo). No obstante, se observa una mayor patrón lineal en los datos que se puede corroborar con la pendiente positiva (aunque pequeña) de la línea de ajuste.

- e) Finalmente, estimen por MCO la ecuación propuesta en el inciso a). A partir de estos resultados, estimen  $\alpha$ ,  $\gamma$ , y sus respectivos errores estándar siguiendo el planteamiento de los incisos b) y c). Presenten en una tabla los parámetros estimados por MCO. ¿Son estos parámetros consistentes con las gráficas que presentaron en el inciso anterior? Discutan brevemente. Adicionalmente, en una segunda tabla presenten los parámetros  $\alpha$  y  $\gamma$  que ustedes estimaron y sus respectivos errores estándar. Interpreten estos últimos. ¿Hay rendimientos crecientes, decrecientes o constantes a escala? ¿Qué tipo de curva de aprendizaje hay en la industria?

■ Solución:

Para estimar por MCO el modelo propuesto en 3, se utilizó el siguiente fragmento de código:

```
1
2 // estimar la ecuación del modelo:
3
4 reg ln_costos ln_producto_acum ln_producto
5 */
6 reg ln_costos ln_producto_acum ln_producto
7
8
9 matlist e(b)
10 matlist e(V)
```

A partir de la estimación de MCO se encuentra que los parámetros estimados son:

$$\hat{\beta}_{MCO} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} -0.0158676 \\ -0.0004461 \\ 0.0053394 \end{pmatrix}$$

Se puede afirmar que los estimadores de MCO son consistentes con las gráficas presentadas en el inciso anterior. En primer lugar, el  $\beta_1$  que es el efecto de  $\log(N_i)$  sobre  $\log(C_i)$  manteniendo todo lo demás constante. Se observa en la gráfica que no hay una relación aparente, lo cual es consistente con el valor de  $\beta_1$  que es cercano a cero. Además, el  $\beta_2$  que es el efecto de  $\log(Y_i)$  sobre  $\log(C_i)$  manteniendo todo lo demás constante, en la gráfica se muestra una relación positiva pero pequeña, lo cual es consistente con el valor del parámetro.

De la estimación se pueden obtener los valores de  $\hat{\alpha}$  y  $\hat{\gamma}$ :

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\gamma} \end{pmatrix} = \begin{pmatrix} \frac{\hat{\beta}_1}{\hat{\beta}_1 + \hat{\beta}_2} \\ \frac{\hat{\beta}_1}{\hat{\beta}_1 + \hat{\beta}_2} \end{pmatrix} = \begin{pmatrix} \frac{-0.0004461}{-0.0004461 + 0.0053394} \\ \frac{-0.0004461}{-0.0004461 + 0.0053394} \end{pmatrix} = \begin{pmatrix} -0.09116547 \\ 204.36107 \end{pmatrix}$$

Finalmente, se observa que hay rendimiento crecientes a escala dado que se obtuvo un valor de  $\gamma = 204.36107 > 1$ , lo que por teoría económica implica que hay rendimiento crecientes a escala. Lo anterior, también significa que hay una *curva de aprendizaje con rendimientos crecientes* porque  $\gamma$  es mayor a 1<sup>5</sup>.

### Tercer ejercicio

Aunque las propiedades asintóticas de los estimadores son de gran utilidad para el ejercicio práctico de la estadística, lo cierto es que muchas veces encontrar resultados teóricos en este ámbito es extremadamente complicado. En estas ocasiones, es usual recurrir a otras herramientas, como lo son los procedimientos de Monte Carlo. En breve, los métodos de Monte Carlo aplicados a la estadística buscan explorar las propiedades de los estimadores (insesgamiento, consistencia, eficiencia, suficiencia, etc.) al observar el comportamiento de los mismos en varias muestras aleatorias simuladas. Este ejercicio los guiará a través de una serie de actividades que les permitirán entender algunos métodos clásicos de simulación de muestras aleatorias para, posteriormente, usarlos para evaluar el comportamiento de uno de los estimadores más importantes que se encontrarán en sus carreras: el estimador de MCO.

Suponga que el Ministerio de Educación ha decidido lanzar un programa de preparación y acompañamiento para la presentación de la prueba ICFES a estudiantes que se encuentren en su último año de bachillerato. En particular, suponga que usted conoce que el puntaje ICFES estandarizado potencial  $Y_i(\cdot)$  obedece el siguiente modelo:

$$Y_i(0) \sim N(0, 1)$$

$$Y_i(1) = Y_i(0) + 3$$

donde  $Y_i(1)$  y  $Y_i(0)$  son los resultados potenciales del puntaje ICFES en caso de participar y de no participar respectivamente, y donde  $Y_i(0)$  tiene una distribución normal estándar.

- a) Escriba el puntaje ICFES observado,  $Y_i$ , en términos de los resultados potenciales y la exposición al tratamiento,  $D_i$ . Suponga que el Ministerio escoge aleatoriamente a los participantes del programa y que todo individuo seleccionado obligatoriamente participa. Proponga un modelo de regresión lineal que le permita estimar el efecto del tratamiento sobre el puntaje, llámese  $\delta$ .

■ Solución:

Para encontrar el puntaje del ICFES observado ( $Y_i$ ), en términos de los resultados potenciales y la exposición al tratamiento,  $D_i$ , se parte del *modelo causal de Rubin*:

$$Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0)$$

$$Y_i = Y_i(0) + (Y_i(1) - Y_i(0)) D_i$$

El *modelo de regresión lineal* que permita estimar el efecto del tratamiento sobre el puntaje ( $\delta$ ) es:

$$Y_i = Y_i(0) + (Y_i(1) - Y_i(0)) D_i$$

$$Y_i = E[Y_i(0)] + (Y_i(1) - Y_i(0)) D_i + E[Y_i(0)] - Y_i(0)$$

$$Y_i = \alpha + \delta D_i + \epsilon_i$$

Donde  $\alpha = E[Y_i(0)]$  y  $\epsilon_i = E[Y_i(0)] - Y_i(0)$

<sup>5</sup>Como se explica en el inciso,  $\gamma$  es el parámetro que caracteriza los retornos a escala de la función de aprendizaje y si la tecnología de producción muestra rendimientos crecientes a escala, los costos unitarios deberían disminuir a medida que aumenta el nivel de producción lo que se refleja claramente con el  $\gamma > 1$

El modelo de regresión lineal  $Y_i = \alpha + \delta D_i + \epsilon_i$  permite estimar el efecto del tratamiento sobre el puntaje ( $\delta$ ), dado que al haber aleatorización se elimina el sesgo de selección (SB), lo que hace que la simple diferencia de medios (SDO), que por definición es equivalente al parámetro que acompaña a la variable de asignación a tratamiento  $D$  que se estima por MCO, sea igual al efecto del tratamiento sobre el puntaje ( $\delta$ )<sup>6</sup>.

- b). Inicialmente, buscaremos validar a través de simulaciones algunas de las propiedades asintóticas del estimador por MCO. En particular, a usted le interesa saber si a medida que aumenta el tamaño de muestra sus estimativos del efecto del programa se tornan más precisos.

■ Solución:

En primer lugar, fue necesario definir la semilla para generar reproducibilidad en el código<sup>7</sup>. Posteriormente, se definieron unas funciones auxiliares para calcular la simple diferencia de medias, la matriz de regresores  $\mathbf{X}$ , el estimador  $\hat{\sigma}$ , el estimador  $\hat{Q}_{XX}^{-1}$  y el estimador para la matriz de varianzas-covarianzas

para errores robustos *Huber-White*  $\hat{avar}(\hat{\beta}_{mco}) = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^N \hat{\epsilon}_i^2 \mathbf{x}_i \mathbf{x}_i' \right) (\mathbf{X}'\mathbf{X})^{-1}$

Posteriormente, se construyó la función *delta*, que es la función más importante del código, dado que es la base para todo lo que se realiza subsecuentemente<sup>8</sup>. La función *delta* está diseñada para calcular entre otras cosas: el *estimador MCO* de  $\delta$  que se obtiene al estimar por MCO la ecuación  $Y_i = \alpha + \delta D_i + \epsilon_i$ , el estimador  $\hat{\sigma}$  y el estimador de la varianza  $\hat{var}(\hat{\delta})$ , dependiendo si se usó la matriz de varianzas-covarianzas tradicional o la matriz de varianzas-covarianzas para errores robustos *Huber-White*.

```
1 library(tidyverse)
2
3 # Definir la semilla para reproducibilidad en los resultados
4 set.seed(12345)
5
6 # Punto 3 ----
7
8 # Ejercicio a. ----
9
10 # Switching regression:
11 ## y_{i} = y_{i}^{0} * D_{i} + (1 - D_{i}) * y_{i}^{1}
12 ## y_{i} = y_{i}^{0} + (y_{i}^{1} - y_{i}^{0}) * D_{i}
13 ## y_{i} = delta * D_{i} + epsilon_{i}
14
15 # Donde el parámetro causal:
16 ## delta = E[y_{i}^{1} | D_{i} = 1] - E[y_{i}^{0} | D_{i} = 0]
17 ## delta = 3 - 0
18 ## delta = 3
19
20 # Ejercicio b. ----
21
22 # Nota: Dado las condiciones para simular y_0 = y_{i}^{0} y y_1 = y_{i}^{1},
23 # uno esperar a que el SDO = 3 y por tanto el delta = 3
24 # Además, dado que el SB = 0 (teniendo en cuenta que y_0 y y_1 son independientes de D)
25 # se tiene que delta representa el impacto causal del programa del ministerio
26
27 # I. ----
28
29 # Definición de parámetros para la simulación
30
31 ## Número de muestras:
32 t = 100
33
34 ## Tamaño de muestras:
35
36 n10 = 10
37 n20 = 20
38 n50 = 50
39 n100 = 100
40 n500 = 500
41 n1000 = 1000
42
43 # -----
44 # Funciones auxiliares:
45 # -----
46
47 # Cálculo del SDO: (SDO := Simple difference of outcomes) (Función auxiliar)
48 SDO_calculo = function(Y, D){
```

<sup>6</sup>Recordar, que como la simple diferencia de medias es  $SDO = \delta + SB$ , si el sesgo de selección (SB) se elimina con la aleatorización, entonces la simple diferencia de medias es igual al efecto causal del tratamiento  $SDO = \delta$

<sup>7</sup>La semilla es una manera en la que podemos replicar código que presente aleatoriedad o que utilice funciones de distribución simuladas, explotando el hecho de que los algoritmos de aleatorización de los computadores son pseudo-aleatorios. De esta forma, es posible generar reproducibilidad en los resultados y obtener los mismos resultados

<sup>8</sup>Esta función es lo suficientemente general que permite hacer las estimaciones de  $\delta$  bajo las diferentes especificaciones del problema. Es decir, permite estimar delta bajo una muestra que distribuye normal, bajo una muestra que distribuye Cauchy y bajo una muestra que distribuye normal pero que presenta heterocedasticidad en los errores. Es la función base, para toda la simulación de Monte Carlo que se realizará subsecuentemente.

```

49 # Por construcci n
50 ## Y = Y_0 Si D = 0
51 ## Y = Y_1 Si D = 1
52 # Nota: Y y D son vectores del mismo tama o
53 # Inicializaci n de contadores
54 count0 = 0
55 count1 = 0
56 # Inicializaci n de la suma para c lculo del SDO
57 suma0 = 0
58 suma1 = 0
59 # Iteraci n a trav s de la muestra
60 for (i in 1:length(Y)){
61   if (D[i] == 0){
62     suma0 = suma0 + Y[i]
63     count0 = count0 + 1
64   }else{
65     suma1 = suma1 + Y[i]
66     count1 = count1 + 1
67   }
68 }
69 SDO = (suma1/count1) -(suma0/count0) # SDO: Simple diferencia de muestras
70 return(SDO)
71 }
72
73 # Construcci n matriz X para un modelo de regresi n lineal con constante: (Funci n
  auxiliar)
74 ols_X = function(...){
75   # ... contiene los regresores necesarios para construir la matriz X
76   regresores = list(...)
77   # Nota: Cada regresor es un vector con el mismo n mero de observaciones n
78   n = length(regresores[[1]]) # No importa que se tome el primer regresor dado que todos los
    regresores tienen el mismo tama o
79   const = rep(1, times = n)
80   X = cbind(const, ...)
81   return(X)
82 }
83
84 # Construcci n del estimador \hat{\sigma}^2\}: (Funci n auxiliar)
85 estimador_sigma = function(y, X, beta_est){
86   # Variables:
87   ## y: es el vector que contiene la variable dependiente
88   ## X: es la matriz que contiene las variables regresoras
89   ## beta_est: es el vector de los par metros estimados por OLS
90   e = y - X %*% beta_est # e es el vector de residuales
91   k = ncol(X) # El n mero de par metros a estimar
92   n = length(y) # El n mero de observaciones en la muestra
93   sigma_est = (t(e) %*% e) / (n - k) # estimador de sigma (varianza del t rmino de error)
94   # El resultado de la funci n es: sigma_est (que es un escalar)
95   return(sigma_est)
96 }
97
98 # Construcci n del estimador para la matriz \hat{Q}_{XX}^{-1}: (Funci n auxiliar)
99 estimador_Q_xx_inv = function(X){
100   n = nrow(X) # donde n es el tama o de la muestra
101   Q = solve((1/n) * (t(X) %*% X)) # Q es el estimador que se est buscando
102   # Q es una matriz de tama o k x k, donde k es en n mero de par metros a estimar
103   # La funci n me retorna a una matriz Q de tama o k x k
104   return(Q)
105 }
106
107 # Construcci n del estimador de la matriz de varianzas y covarianzas
  # Huber-White: (Funci n auxiliar)
108 Huber_White = function(y, X, beta_est){
109   # Variables:
110   ## y: es el vector que contiene la variable dependiente
111   ## X: es la matriz que contiene las variables regresoras
112   ## beta_est: es el vector de los par metros estimados por OLS
113   e = y - X %*% beta_est # e es el vector de residuales
114   k = ncol(X) # El n mero de par metros a estimar
115   n = length(y) # El n mero de observaciones en la muestra
116   # Construcci n de (X^{'}X)^{-1}
117   mat_X = solve(t(X) %*% X)
118   # Matriz intermedia de los errores robustos White (que se encuentra por medio de una suma)
119   # Inicializo la matriz como una matriz de ceros
120   mat_intermedia = matrix(0, nrow = k, ncol = k)
121   # El for se dise a para llenar la matriz mat_intermedia
122   for (i in 1:n){
123     mat_intermedia = mat_intermedia + (e[i])^2 * (X[i, ] %*% t(X[i, ]))
124   }
125   # mat_final me dal estimador de la matriz de varianzas y covarianzas de Huber-White
126   mat_final = mat_X %*% mat_intermedia %*% mat_X
127   return(mat_final)
128 }
129
130
131 # y = rnorm(1000)
132 # D = rnorm(1000, mean = 2)
133 # beta = c(1, 2)
134 # X = ols_X(D)
135 #
136 # Huber_White(y, X, beta)
137
138

```



```

139 # l estimador para la matriz \Hat{Q}_{XX}^{-1}: (Funci n auxiliar)
140
141 ##-----
142 # delta: Es la funci n principal (m s importante) de todo el c digo.
143 # Con base en la funci n delta es que se deduce todo lo dem s
144 # en el c digo.
145 # Es una funci n bastante general que contempla todas las situaciones
146 # que puedan surgir en el c digo. En ese caso, contempla la simulaci n
147 # de la Cauchy y de una muestra heteroced stico donde la volatilidad
148 # de los tratados es diferente a la volatilidad de los no tratados
149 # Nota: Revisar toda el enunciado del ejercicio primero para entender mejor
150 # la l gica de la funci n
151 ##-----
152
153 # Funci n para estimar el efecto tratamiento (delta) en una base de datos
154 delta = function(t, n, distro, hetero = F, varianza_y0 = 1, varianza_y1 = 1, robustos = F){
155   # Defino la semilla para reproducibilidad del resultado
156   set.seed(5678)
157   # Definici n de variables:
158   ## t: n mero de muestras a simular
159   ## n: tama o de las muestras
160   ## distro: tipo de distribuci n con el que se va a simular Y_{i}^{0}
161   ## Nota: Si distro == "normal", entonces los par metros opcionales
162   ## empiezan a tomar relevancia
163   ## hetero: Para saber si la volatilidad entre tratados y no tratados es diferente
164   ## varianza_y0: El valor de la varianza de la distribuci n normal con la que se simula y_{i}^{0}
165   ## varianza_y1: El valor de la varianza de la distribuci n normal con la que se simula y_{i}^{1}
166   ## robustos: Solo se activa si se escoge la opci n hetero == T y adem s robustos == T
167   ## Permite estimar errores robustos usando la matriz de varianzas y covarianzas
168   # Consideraciones adicionales:
169   # Condici n l gica para saber si y_0 sigue una distribuci n normal est ndar o una
170   # distribuci n Cauchy est ndar
171   if (distro == "normal"){
172     # Hay dos posibles casos en caso de que la distribuci n sea normal o no:
173     # 1. Homoced sticidad entre tratados y no tratados
174     ## En caso que la volatilidad de y_{i}^{0} y y_{i}^{1} sea la misma hetero == F
175     # 2. Heterocedastidad entre tratados y no tratados (diferente volatilidad dependiendo si
176     ## fueron tratados o no)
177     ## En caso que la volatilidad de y_{i}^{0} y y_{i}^{1} sea diferente hetero == T
178     if (hetero == F){
179       # df: Dataframe que almacena el SDO y el delta_estimado por OLS
180       # El par metro delta es el para etro asociado a la asignaci n a tratamiento (D)
181       df = data.frame(SDO = double(), delta_est = double(), sigma = double(), var_tlc =
182         double(), var_ols = double())
183       # Meter un for para generar las t diferentes muestras
184       for (muestra in 1:t){
185         # Simulaci n de los outcomes potenciales y de la variable trataminto
186         y_0 = rnorm(n, mean = 0, sd = 1) # Simulaci n del outcome potencial de ausencia de
187         tratamiento
188         y_1 = y_0 + 3 # Simulaci n del outcome potencial de presencia de tratamiento
189         D = rbinom(n, 1, prob = 0.3) # Simulaci n de una variable Bernoulli con una
190         probabilidad de xito de 0.3
191         # Modelo causal de Rubin
192         y = y_0 + (y_1 - y_0) * D # y es un vector Nx1, donde n es el n mero de
193         observaciones (tama o de muestra)
194         # Corregir (Debo calcular es el SDO)
195         SDO = SDO_calculo(y, D) # C lculo del SDO utilizando la funci n SDO_calculo
196         # Estimaci n del par metro delta mediante una regresi n lineal con constante
197         # X es una matriz NxK, donde n es el n mero de observaciones y K el n mero de
198         para etros
199         X = ols_X(D) # D es la variable tratamiento
200         beta_ols = solve(t(X) %*% X) %*% t(X) %*% y # beta_ols es un vector de Kx1, donde K
201         es el n mero de par metros
202         # Ahora, se calcularon 3 par metros adicionales:
203         # el estimador de sigma, el estimador de var(sqrt(n) (delta_{OLS} - delta) y el
204         # estimador de var(delta_{OLS}))
205         ## estimador de sigma:
206         sigma = estimador_sigma(y, X, beta_ols)
207         ## estimador de \Hat{Q}_{XX}^{-1}:
208         Q_mat = estimador_Q_xx_inv(X)
209         ## estimador de var_tlc = var(sqrt(n) (delta_{OLS} - delta))
210         var_tlc = sigma * Q_mat[[2,2]]
211         ## estimador de var_ols = var(delta_{OLS})
212         n = nrow(X)
213         var_ols = 1/n * sigma * Q_mat[[2,2]]
214         # Dataframe que contiene todos los par metros de inter s
215         df[muestra, ] = c(SDO, beta_ols[[2,1]], sigma, var_tlc, var_ols) # Extraigo el
216         segundo par metro que es el par metro delta
217       }
218     }else{
219       # df: Dataframe que almacena el SDO y el delta_estimado por OLS
220       # El par metro delta es el para etro asociado a la asignaci n a tratamiento (D)
221       df = data.frame(SDO = double(), delta_est = double(), sigma = double(), var_ols =
222         double())
223       # Meter un for para generar las t diferentes muestras
224       for (muestra in 1:t){
225         # Simulaci n de los outcomes potenciales y de la variable trataminto
226         D = rbinom(n, 1, prob = 0.3) # Simulaci n de una variable Bernoulli con una
227         probabilidad de xito de 0.3

```

```

216     y_0 = rnorm(n, mean = 0, sd = sqrt(varianza_y0)) # Simulaci n del outcome potencial
217     y_1 = rnorm(n, mean = 3, sd = sqrt(varianza_y1)) # Simulaci n del outcome potencial
218     de presencia de tratamiento
219     # Modelo causal de Rubin
220     y = y_0 + (y_1 - y_0) * D # y es un vector Nx1, donde n es el n mero de
221     observaciones (tama o de muestra)
222     # Corregir (Debo calcular es el SDO)
223     SDO = SDO_calculo(y, D) # C lculo del SDO utilizando la funci n SDO_calculo
224     # Estimaci n del par metro delta mediante una regresi n lineal con constante
225     # X es una matriz NxK, donde n es el n mero de observaciones y K el n mero de
226     para etros
227     X = ols_X(D) # D es la variable tratamiento
228     beta_ols = solve(t(X) %*% X) %*% t(X) %*% y # beta_ols es un vector de Kx1, donde K
229     es el n mero de par metos
230     if (robustos == F){
231       # Ahora, se calcularon 3 par metros adicionales:
232       # el estimador de sigma, el estimador de var(sqrt(n) (delta_{OLS} - delta) y el
233       estimador de var(delta_{OLS}))
234       ## estimador de sigma:
235       sigma = estimador_sigma(y, X, beta_ols)
236       ## estimador de \Hat{Q}_{XX}^{-1}:
237       Q_mat = estimador_Q_xx_inv(X)
238       ## estimador de var_ols = var(delta_{OLS})
239       n = nrow(X)
240       var_ols = 1/n * sigma * Q_mat[[2,2]]
241       # Dataframe que contiene todos los par metros de inter s
242       df[muestra, ] = c(SDO, beta_ols[[2,1]], sigma, var_ols) # Extraigo el segundo
243       par metro que es el par metro delta
244     }else{
245       # Sigma, es provisional, luego se retira porque no tiene sentido para un estimador
246       con matriz de varianzas y covarianzas
247       # con errores robustos Huber-White (Dado la heterocedasticidad de los errores)
248       sigma = 0
249       # C lculo errores robustos Huber-White (me genera una matriz k x k, donde k es el
250       n mero de par metros)
251       var_ols = Huber_White(y, X, beta_ols)[[2,2]] # Extraigo la componente 2 de la
252       matriz de varianzas y covarianzas de Huber-White
253       df[muestra, ] = c(SDO, beta_ols[[2,1]], sigma, var_ols) # Extraigo el segundo
254       par metro que es el par metro delta
255     }
256   }
257   # Condicional final para eliminar la columna sigma si robustos == T
258   if (robustos == T){
259     df = df %>%
260       select(SDO, delta_est, var_ols)
261   }
262 }
263 }else if (distro == "cauchy"){
264   # df: Dataframe que almacena el SDO y el delta_estimado por OLS
265   # El par metro delta es el para etro asociado a la asignaci n a tratamiento (D)
266   df = data.frame(SDO = double(), delta_est = double())
267   for (muestra in 1:t){
268     # Simulaci n de los outcomes potenciales y de la variable trataminto
269     y_0 = rcauchy(n, location = 0, scale = 1) # Simulaci n del outcome potencial de
270     ausencia de tratamiento
271     y_1 = y_0 + 3 # Simulaci n del outcome potencial de presencia de tratamiento
272     D = rbinom(n, 1, prob = 0.3) # Simulaci n de una variable Bernoulli con una
273     probabilidad de xito de 0.3
274     # Modelo causal de Rubin
275     y = y_0 + (y_1 - y_0) * D # y es un vector Nx1, donde n es el n mero de
276     observaciones (tama o de muestra)
277     # Corregir (Debo calcular es el SDO)
278     SDO = SDO_calculo(y, D) # C lculo del SDO utilizando la funci n SDO_calculo
279     # Estimaci n del par metro delta mediante una regresi n lineal con constante
280     # X es una matriz NxK, donde n es el n mero de observaciones y K el n mero de
281     para etros
282     X = ols_X(D) # D es la variable tratamiento
283     beta_ols = solve(t(X) %*% X) %*% t(X) %*% y # beta_ols es un vector de Kx1, donde K es
284     el n mero de par metos
285     df[muestra, ] = c(SDO, beta_ols[[2,1]]) # Extraigo el segundo par metro que es el
286     par metro delta
287   }
288 }
289 # La funci n delta retorna un df con el c lculo de la SDO
290 # y el par metro estimado delta, que proviene de una switching regression
291 return(df)
292 }

```

Luego de tener listas las *funciones auxiliares* y la *funci3n principal delta*, se procede a realizar las simulaciones.

- I. Simule 100 muestras  $\{(Y_i(0), Y_i(1), D_i)_{i=1}^n\}$  de tama1o  $n = 10, 20, \dots, 1000$  donde  $D_i \sim \text{Bernoulli}(0.3)$ . Para cada muestra, estime el efecto de la pol3tica del Ministerio sobre el puntaje ICFES y almacene su estimado.

- Soluci3n:

El siguiente fragmento de c3digo permite realizar las simulaciones. En cada simulaci3n, se simulan

100 muestras, en donde lo que cambian son los tamaños empleados<sup>9</sup>.

```
1
2 # Simulaci n
3 # delta(t, n10, distro = "normal") # Existe el riesgo de que haya
4 # multicolinealidad perfecta cuando se usa n10 = 10
5 # Puede generar un vector D = rep(0, times = 10)
6
7 delta20 = delta(t, n20, distro = "normal")
8 delta50 = delta(t, n50, distro = "normal")
9 delta100 = delta(t, n100, distro = "normal")
10 delta500 = delta(t, n500, distro = "normal")
11 delta1000 = delta(t, n1000, distro = "normal")
```

II. Obtenga la media y la varianza muestral de los estimadores almacenados para los distintos tamaños de muestra.

- Solución:

Para encontrar la media y la varianza muestral de los estimadores se empleo el siguiente fragmento de código:

```
1
2 # n_vector es un vector que almacena los diferentes tama os de muestra
3 n_vector = seq(from = 20, to = 1000, by = 10)
4
5 # Funci n que genera un dataframe con el tama o de muestra,
6 # la media y la varianza del estimador de delta para cada
7 # tama o diferente de muestra
8
9 media_varianza_delta = function(t, n_vector, distro, hetero = F, varianza_y0 = 1,
10   varianza_y1 = 1, robustos = F){
11   # Variables:
12   ## t: N mero de muestras que se van a generar por cada tama o de muestra
13   ## n_vector: Variable que almacena los diferentes tama os de muestras
14   ## distro, hetero, varianza_y0 y varianza_y1 son par metros definidos para la
15   funci n delta
16   # df: DataFrame que almacena el tama o de muestra,
17   # la media y la varianza del estimador de delta para cada tama o de muestra
18   df = data.frame(tama o = double(), media = double(), varianza = double())
19   # Llamo a la funci n delta para cada tama o diferente de muestra
20   for (i in 1:length(n_vector)){
21     n_muestra = n_vector[i]
22     delta_n = delta(t, n_muestra, distro, hetero, varianza_y0, varianza_y1, robustos)$
23     delta_est # delta_n es el vector de deltas por cada tama o de muestra
24     df[i,] = c(n_muestra, mean(delta_n), var(delta_n))
25   }
26   return(df)
27 }
28
29 # Dataframe con el tama o de muestra y la media y varianza del
30 # estimador de delta para los diferentes tama os de muestra
31 media_varianza = media_varianza_delta(100, n_vector, distro = "normal"); glimpse(media_
32   varianza)
```

Lo que se puede observar es el estimador es insesgado<sup>10</sup> y consistente<sup>11</sup>. Igualmente, se observa que la varianza muestral de  $\delta$  cae a medida que aumenta el tamaño de las muestras.

III. Haga las siguientes dos gráficas: 1) Grafique el promedio muestral de los estimadores contra el tamaño de muestra. 2) Grafique las varianzas muestrales contra el tamaño de muestra correspondientes. ¿Qué puede decir acerca del comportamiento de las estimaciones a medida que aumenta el tamaño de muestra? ¿Qué propiedades del estimador de MCO se ven reflejadas en el ejercicio?

- Solución:

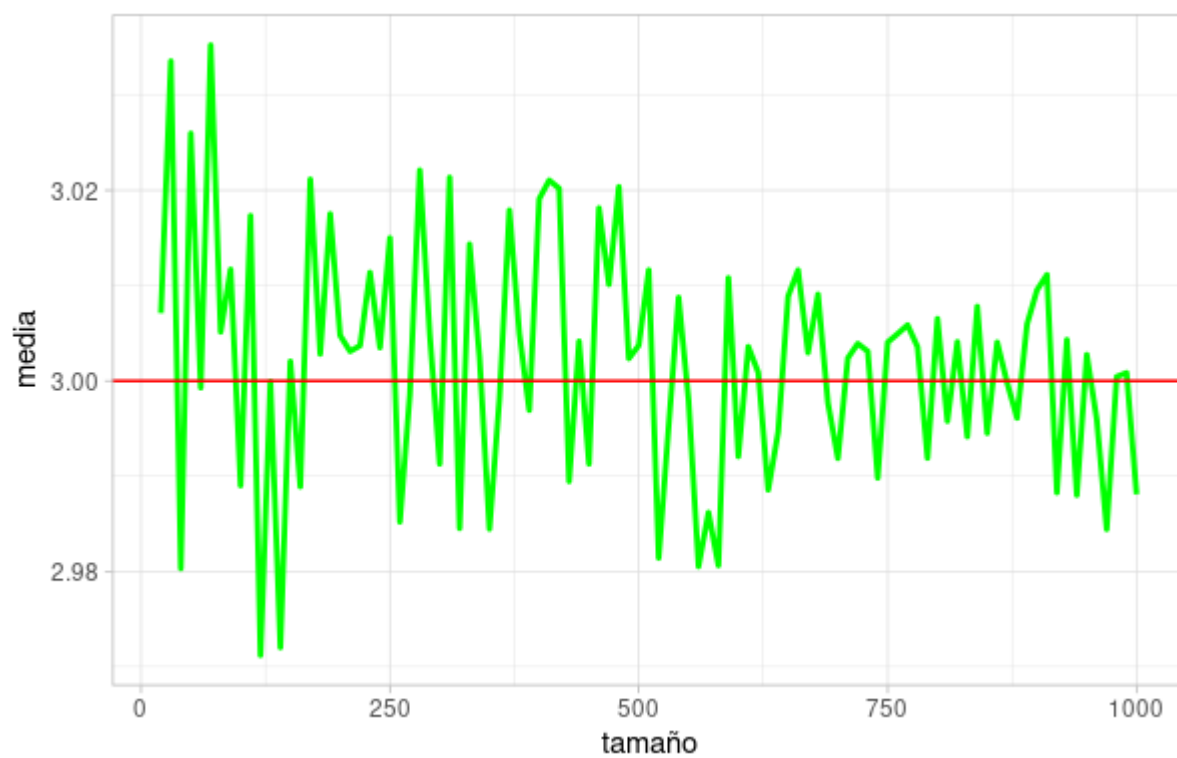
La gráfica de *promedio muestral* de los estimadores contra el tamaño de muestra:

<sup>9</sup>Se utilizan tamaños de muestras de: 20, 50, 100, 500 y 1000 observaciones. No se utilizan tamaños de muestras de 10 observaciones porque al tener la asignación al tratamiento un parámetro de probabilidad de éxito de  $\theta = 0.3$ , hace factible generar un vector de la variable tratamiento donde ninguno de las 10 observaciones reciba tratamiento, lo que hace que no haya variabilidad en el regresor y genere multicolinealidad perfecta con la variable que modela la constante del modelo.

<sup>10</sup>Con muestra prácticamente de tamaño 20 ya converge al verdadero valor = 3

<sup>11</sup>Con muestra muy grandes los estimadores se aproximan muy bien a = 3

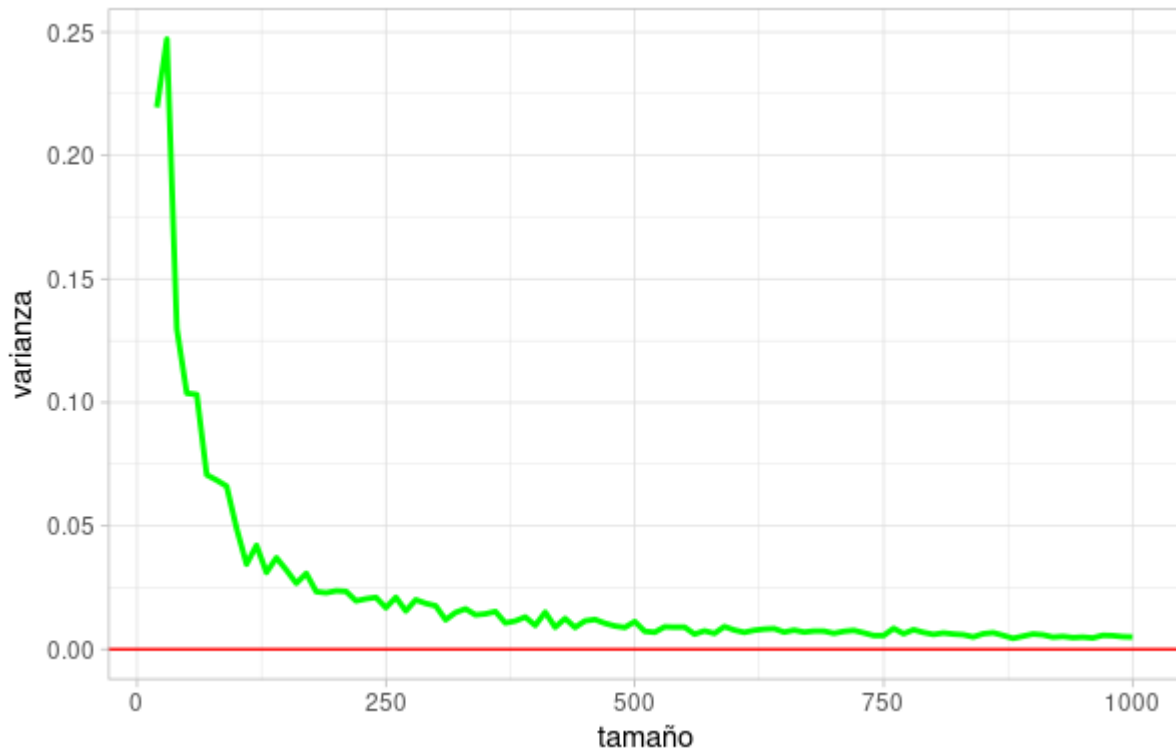
Figura 3: Gráfica promedio muestral de los estimadores contra el tamaño de muestra.



Variables incluidas: media: promedio muestral de los estimadores MCO de  $\delta$  y tamaño: tamaño de muestra

La gráfica de la varianza muestral de los estimadores contra el tamaño de muestra:

Figura 4: Gráfica varianza muestral de los estimadores MCO de  $\delta$  contra el tamaño de muestra.



Variables incluidas: varianza: varianza muestral de los estimadores y tamaño: tamaño de muestra

Para construir las dos gráficas: 1) el promedio muestral de los estimadores contra el tamaño de muestra y 2) las varianzas muestrales contra el tamaño de muestra correspondiente se empleó el siguiente fragmento de código:

```

1  grafica_propiedades = function(df, variable_y, y_intercepto = 0){
2    variable_y = ensym(variable_y)
3    graph = df %>%
4      ggplot(aes(x = tamaño, y = !!variable_y)) +
5      geom_line(color = "green", size = 1) +
6      geom_hline(yintercept = y_intercepto, color = "red") +
7      theme_light()
8    return(graph)
9  }
10 }
11
12 grafica_propiedades(media_varianza, media, y_intercepto = 3)
13 grafica_propiedades(media_varianza, varianza)

```

Ahora bien, para contextualizar, se sabe que por las instrucciones dadas, para cada tamaño de muestra empleada, se simuló 100 muestras. Para la generación de las gráficas se construyó una especie de *grid* o rejilla que permitiera generar muestras que iban de tamaño 20 a tamaño 1000, aumentando el tamaño de la muestra de 10 en 10 observaciones, pero siempre generando 100 muestras por cada tamaño de muestra<sup>12</sup>. Con esta metodología, y como se está manteniendo siempre el número de muestras simuladas por tamaño de muestra de  $t = 100$ , es posible visualizar cuáles son los efectos de aumentar el tamaño de las muestras sobre los estimadores MCO del parámetro de interés  $\delta$ . Como se puede observar de 3, a medida que  $n$  aumenta se puede observar por la serie gráfica que la media muestral de dichos estimadores oscila alrededor del parámetro poblacional conocido  $\delta = 3$ . Aún más interesante, es observar que desde  $n = 20$  la media de dichos estimadores es muy cercana a 3 lo que muestra la propiedad de insesgadez del estimador MCO y que prácticamente a lo largo de la serie se observa que dichas medias muestrales del estimador rara vez exceden los valores de 3.03 o caen por debajo de 2.97, lo que muestra nuevamente la alta precisión del estimador MCO para todos los tamaños de muestra. De igual forma, de la 3 se observa que a medida que el tamaño de la muestra aumenta la volatilidad de la media muestral del estimador MCO disminuye y cada vez dicha media se acerca más y más a  $\delta = 3$ , lo que permite ver

<sup>12</sup>Es importante notar que en una simulación de Monte Carlo hay dos parámetros clave, un parámetro  $t$  que denota el número de muestras que se está simulando y un parámetro  $n$  que denota el tamaño de cada muestra

la propiedad asintótica del estimador MCO de consistencia  $\text{plim } \delta_{MCO} = \delta = 3$ . Finalmente, de la figura 4, claramente se observa que la varianza muestral del estimador MCO  $\delta_{MCO}$  de  $\delta$  disminuye a medida que aumenta el tamaño de muestra, donde claramente se ve una varianza muestral de hasta 0.25 en las primeras muestras, que rápidamente a medida que aumenta el tamaño de las muestras (manteniendo el número de muestras simuladas por tamaño de muestra constante en  $t = 100$ ) cae a 0.05 en las muestras de tamaño  $n = 100$  y ya por muestras con tamaño por encima a 250 observaciones se observa prácticamente que la varianza del estimador MCO de  $\delta$  es casi indistinguible de cero, lo cuál claramente muestra que el estimador MCO es más eficiente a medida que aumenta el tamaño de muestra y que la distribución asintótica del estimador MCO de  $\delta$  tiende a una distribución normal degenerada centrada en el parámetro poblacional  $\delta = 3$ .

Uno de los propósitos principales de los modelos de regresión lineal es poder estudiar las relaciones que existen entre la variable dependiente y las variables independientes. Más precisamente, nos interesa dilucidar el tamaño y el signo de dicha relación. No obstante, en la práctica, esto es difícil de establecer puesto que desconocemos el proceso generador de los datos y en general, contamos con una sola muestra para nuestra estimación. Así las cosas, si nuestro objetivo es hacer inferencia, debemos preguntarnos qué tan precisas son nuestras estimaciones, esto es, qué tan dependientes son de la muestra particular que tenemos. Para lograrlo, es útil entender a los estimadores como variables aleatorias, que dependen de una muestra, pero que tienen una distribución definida. En particular, es de nuestro interés dicha distribución cuando el tamaño de muestra es grande, pues conocer la distribución exacta para muestras reducidas puede ser complicado.

c) En este inciso vamos a aproximar la distribución asintótica del estimador de MCO. Para ello, realicen las siguientes instrucciones:

- I. Siguiendo los pasos expuestos en el punto b), simule 1000 muestras  $\{(Y_i(0), Y_i(1), D_i)\}_{i=1}^n$  de tamaño  $n = 10, 20, 100, 1000$  donde  $D_i \sim \text{Bernoulli}(0.3)$ . Para cada muestra, estime el efecto del programa del Ministerio y calcule y almacene

$$a_{k,n} = \sqrt{n}(\hat{\delta}_{k,n} - \delta)$$

donde  $k$  indexa las muestras de un determinado tamaño.

- Solución:

Dado que se va a estudiar las *propiedades asintóticas del estimador MCO* es necesario recurrir a teoría y teoremas de teoría asintótica para realizar dicho análisis. Lo primero, es recurrir al *teorema del limite central* para encontrar la distribución asintótica del estimador  $\delta_{MCO}$ . Se sabe que:

$$\sqrt{N}(\hat{\beta}_{MCO} - \beta) \xrightarrow{d} \mathcal{N}(0, \sigma^2 Q_{XX}^{-1}) \quad (6)$$

Dicha ecuación, implica que la distribución asintótica de  $\hat{\beta}_{MCO}$  sea:

$$\hat{\beta}_{MCO} \xrightarrow{d} \mathcal{N}(\beta, \sigma^2 Q_{XX}^{-1}) \quad (7)$$

Ahora bien, partiendo de la ecuación 7 se puede obtener la distribución asintótica del estimador  $\hat{\delta}_{MCO}$ .

Como lo muestra la gráfica 5, claramente, el estimador  $\hat{\delta}_{MCO}$  está centrado en 3 y su varianza tiende a disminuir a medida que aumente  $N$  lo cuál hace pensar que la distribución asintóticamente converge a una normal degenerada centrada en  $\delta = 3$ , que es el parámetro poblacional.

Ahora bien, como lo que se está solicitando es:

$$a_{k,n} = \sqrt{n}(\hat{\delta}_{k,n} - \delta)$$

el siguiente código calcula dicho  $a_{k,n}$  y almacena sus valores:

```

1
2 # Ejercicio c. ----
3
4 # I. ----
5
6 # Definición de parámetros para la simulación
7
8 ## Número de muestras:
9 t2 = 1000
10
11 ## Tamaño de muestras:
12
13 # No se puede trabajar con n10 = 10 porque no satisface la condición de rango

```

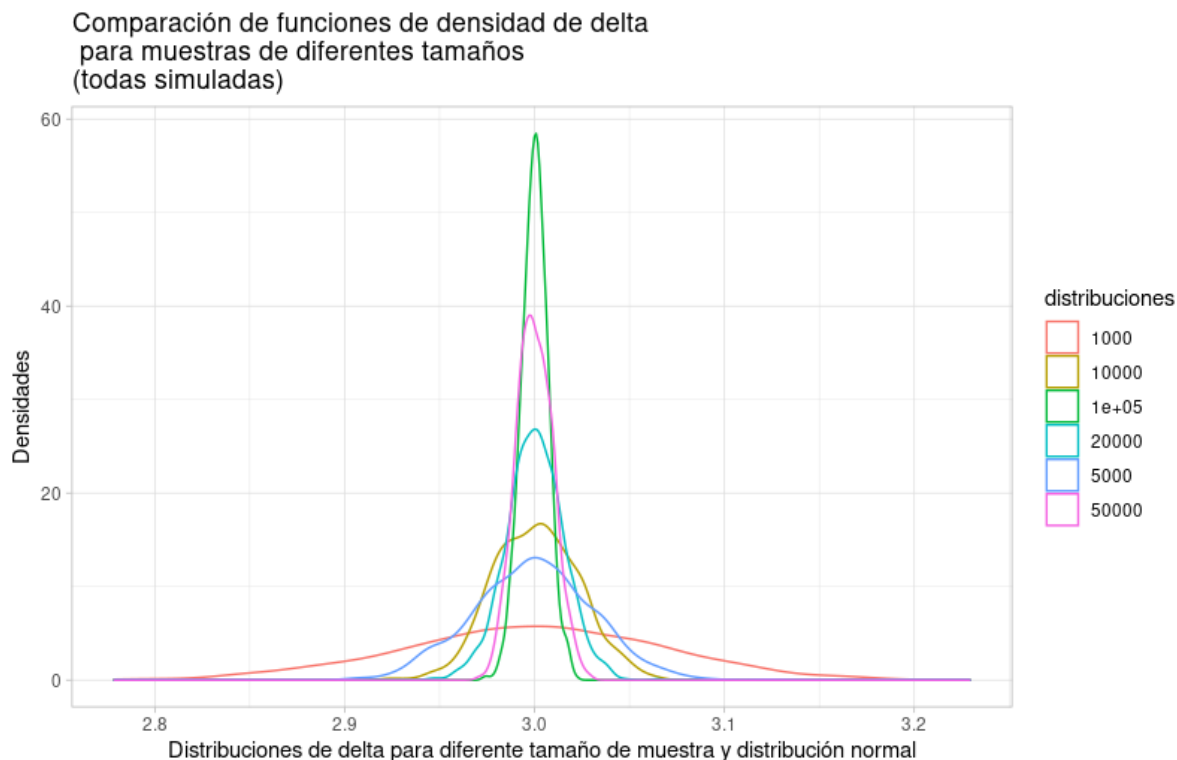


Figura 5: Comportamiento de la distribución asintótica del estimador  $\delta_{MCO}$  para una muestra que distribuye normal a medida que el tamaño de la muestra aumenta

```

14 # n10 = 10 Genera problemas de singularidad cuando se generan muchas muestras
15 #     Porque por chance, se genera una muestra donde D = rep(0, times = 10)
16 n20 = 20
17 n100 = 100
18 n1000 = 1000
19 n5000 = 5000
20 n10000 = 10000
21 n20000 = 20000
22 n50000 = 50000
23 n100000 = 100000
24 n_vector2 = c(n20, n100, n1000, n5000, n10000, n20000, n50000, n100000)
25
26 # Función que calcula a = sqrt(n) * (delta_est - truth_delta) para cada tamaño de
27 # muestra n
28 calculo_a = function(t, n, distro, truth_delta){
29   # Variables:
30   ## t: número de muestras a simular
31   ## n: tamaño de la muestra
32   ## truth_delta: verdadero valor del parámetro delta (valor poblacional del parámetro)
33   # df_estimados es el dataframe que contiene el delta estimado por SDO
34   # y por medio de regresión
35   df_estimados = delta(t, n, distro)
36   # df_a es el dataframe que contiene
37   df_a = df_estimados %>%
38     mutate(a = sqrt(n) * (delta_est - truth_delta), tamaño = rep(n, times = t))
39   return(df_a)
40 }
41
42 # Defino una función para crear un dataframe que se va a utilizar
43 # Para construir la gráfica multipanel.
44 base_para_graficas = function(t, n_vector, distro, truth_delta){
45   for (i in 1:length(n_vector)){
46     # Si i == 1 significa que estamos en la primera iteración
47     if (i == 1){
48       df_a_total = calculo_a(t, n_vector[i], distro, truth_delta = 3)
49     }else{
50       df_prov = calculo_a(t, n_vector[i], distro, truth_delta = 3)
51       df_a_total = bind_rows(df_a_total, df_prov)
52     }
53   }
54   # Retorna un dataframe con los diferentes valores de a para cada tamaño de muestra
55   listo
56   # para construir la gráfica multipanel
57   return(df_a_total)
58 }
59
60 # La base multipanel me tiene lista los diferentes a
61 # para cada tamaño de muestra
62 base_multipanel = base_para_graficas(t2, n_vector2, distro = "normal", truth_delta = 3)

```

II. Grafique las densidades estimadas de los  $a_{k,n}$  para cada  $n$ . ¿Qué puede apreciar a medida que aumenta  $n$ ? ¿A qué se debe este resultado?

• Solución:

Las gráficas de las densidades estimadas de los  $a_{k,n}$  para cada tamaño de muestra  $n$  están dadas por<sup>13</sup> la figura 6.

Comparación de funciones de densidad de  $a$  para muestras de diferentes tamaños (todas simuladas)

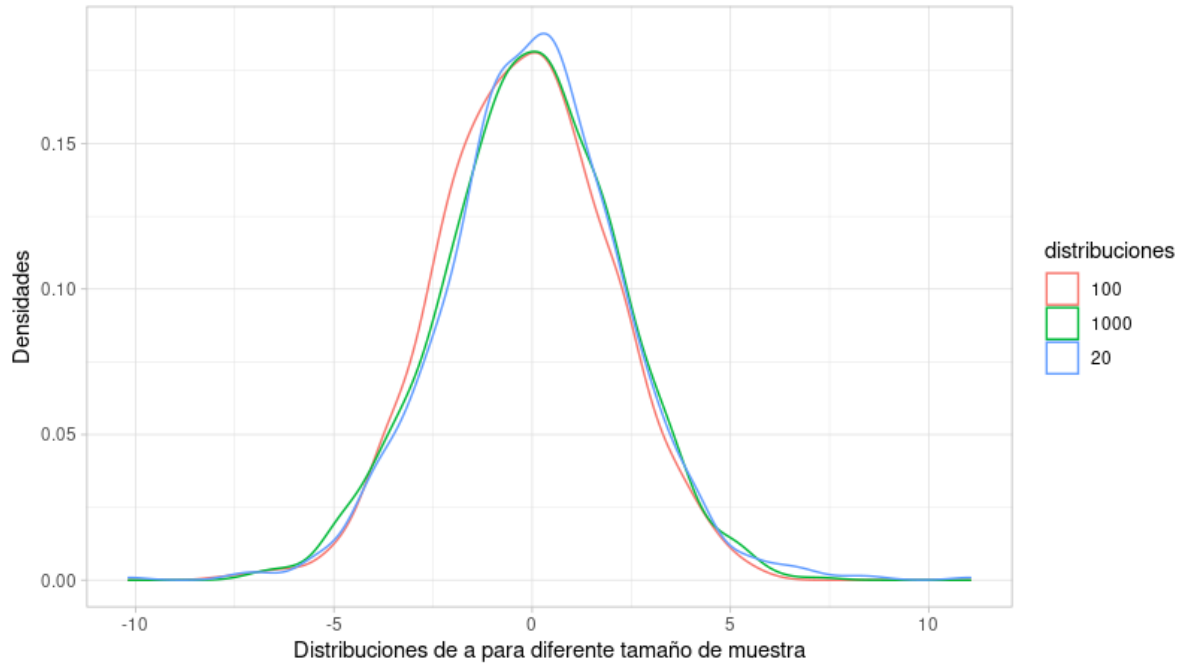


Figura 6: Distribuciones asintóticas de  $a_{k,n}$  en donde claramente se satisface el teorema de limite central para una muestra que distribuye normal.

Lo que se puede apreciar de la figura 6 es que claramente cuando se tiene muestras normales, el *teorema del limite central* se satisface para la media muestral del estimador  $\delta_{MCO}$ . Como se puede observa para todo tamaño de muestra, se observa claramente una distribución normal centrada en cero y que a medida que  $n$  aumenta se observa una convergencia de la distribución. Nuevamente, dicha convergencia se debe principalemnte a que la media muestral del estimador  $\delta_{MCO}$  satisface el *teorema del limite central*.

El código para generar la gráfica fue el siguiente:

```
1
2 # Defino una función para crear la gráfica multipanel con el histograma y la función
  de densidad
3 # para cada tamaño de muestra n
4 histogram_grid = function(df, titulo, x_lab, y_lab, num_bins = 30, y_upper_limit){
5   histog_grid = df %>%
6     ggplot(aes(a)) +
7     scale_y_continuous(limits = c(0, y_upper_limit)) +
8     geom_histogram(aes(y = ..density..), color = "black", bins = num_bins) +
9     geom_density(color = "green") +
10    facet_grid(cols = vars(tamaño)) +
11    ggtitle(titulo) +
12    xlab(x_lab) +
13    ylab(y_lab) +
14    theme_light()
15 }
16
17 grafica_multi = histogram_grid(df = base_multipanel, titulo = "Gráfica multipanel para
  a con diferentes tamaños de muestra", x_lab = "Tamaño de muestras (n)", y_lab = ""
  , num_bins = 30, y_upper_limit = 0.2); grafica_multi
18
19 # Gráfica que compara directamente las densidades
20 comparacion_densidades_grafica = function(t, df, var){
21   # Variables:
22   ## t: número de muestras simuladas por tamaño de muestra
```

<sup>13</sup>Por las instrucciones del taller, se emplearon muestras de tamaño 20, 100, 1000. Las muestras de tamaño  $n = 10$  no se emplearon por los problemas de multicolinealidad comentados anteriormente



```

23 ## df: df multipanel que tiene las simulaciones de las diferentes muestra por tama o
    de muestra
24 ## var: variable simulada de la cual se va a generar las funciones de densidad
25 # Nota:
26 # Modifico el dataframe multipanel para seleccionar solo
27 # las variables a y tama o que permiten hacer la gr fica de densidades
28 # df_mod es un dataframe modificado del dataframe multipanel
29 df_mod = df %>%
30   rename(valores_sim = {{ var }}, distribuciones = tama o) %>%
31   mutate(distribuciones = as.character(distribuciones)) %>%
32   select(valores_sim, distribuciones)
33 # Simulo de una distribuci n normal est ndar porque quiero probar
34 # si se cumple o no el teorema del l mite central
35 df_normal = data.frame(valores_sim = rnorm(t), distribuciones = rep("Normal est ndar"
    , times = t))
36 # base_grafica ya es la base de datos lista para realizar las gr ficas de
37 # las densidades de las a y de una normal est ndar
38 base_grafica = df_mod
39 # base_grafica = bind_rows(df_mod, df_normal)
40 # Gr fica de densidades
41 density_comparacion = base_grafica %>%
42   ggplot(aes(x = valores_sim, color = distribuciones)) +
43   geom_density() +
44   theme_light() +
45   ggtitle("Comparaci n de funciones de densidad de a\n para muestras de diferentes
    tama os \n(todas simuladas)") +
46   ylab("Densidades") +
47   xlab("Distribuciones de a para diferente tama o de muestra")
48   return(density_comparacion)
49 }
50
51 # Visualizaci n de la base de datos con las variables simuladas
52 # delta_est y a
53 glimpse(base_multipanel)
54
55 base_filtrada = base_multipanel %>%
56   filter(100 < tama o)
57
58 base_filtrada_a = base_multipanel %>%
59   filter(10 < tama o) %>%
60   filter(tama o < 2000)
61
62 # Gr ficas de las funciones de densidad para la simulaci n de los delta estimados
63 # y de los a, para el caso de un  $Y_{i|0} = \text{norm}(0, 1)$ 14
64 grafica_a = compracion_densidades_grafica(t = t2, df = base_filtrada_a, var = a);
    grafica_a
65 grafica_delta_est = compracion_densidades_grafica(t = t2, df = base_filtrada, var =
    delta_est); grafica_delta_est

```

Suponga ahora que usted sabe que existen factores ajenos a las políticas del Ministerio que inciden sobre los resultados potenciales de los alumnos. Por ejemplo, sabe que existen alumnos que, independientemente de si son o no seleccionados, igual se inscribirían en cursos de preparación para presentar el ICFES. Similarmente, hay alumnos que por condiciones adversas (enfermedades, salones menos adecuados para la presentación del examen) exhiben un desempeño muy inferior a lo esperado. Estos factores producen una mayor volatilidad en los resultados potenciales de los alumnos.

- d) En este sentido, usted sabe que es más plausible considerar que los resultados potenciales en realidad obedecen la ley  $Y_i(0) \sim \text{Cauchy}$  (**Distribución Cauchy estándar**) en vez de una normal estándar. Esto es, en ocasiones usted observa datos atípicos que no parecen corresponder con el patrón general de la muestra. Repita los incisos I y II del c) bajo estas condiciones. ¿Por qué en este caso no se satisface el Teorema del Límite Central? ¿Qué consecuencias tendría esto para la inferencia de los parámetros de MCO si usáramos los procedimientos estadísticos usuales en este caso?

■ Solución:

Ahora bien, para la solución de este ejercicio se volvieron a repetir los incisos I y II del punto c) pero asumiendo que el outcome potencial en ausencia de tratamiento distribuye como una  $Y(0) \text{ Cauchy}(0, 1)$ <sup>14</sup>. La idea de utilizar una distribución Cauchy, es que permite observar datos atípicos que no parecen corresponder con el patrón general de la muestra.

Lo primero, fue simular el estimador de MCO  $\hat{MCO}$  bajo la nueva muestra que distribuye Cauchy.

La gráfica 7 muestra la distribución asintótica del estimador  $\hat{MCO}$ . Lo primero que es evidente es que es que se observa que dicho estimador *distribuye como una Cauchy*. Esto es particularmente cierto, si se

<sup>14</sup>Una distribución  $\text{Cauchy}(0, 1)$ , es una distribución Cauchy estándar con parámetro de localización 0 y parámetro de escala 1 (que le da la forma a la distribución)



Figura 7: Comportamiento de la distribución asintótica del estimador  $\delta_{MCO}$  para una muestra que distribuye normal a medida que el tamaño de la muestra aumenta

tiene en cuenta la cantidad muy alta de observaciones que se encuentran alrededor del parámetro de localización 0, pero también por que existen valores atípicos muy extremos<sup>15</sup>.

Ahora bien, las gráficas de las densidades estimadas de los  $a_{k,n}$  para cada tamaño de muestra  $n$ , donde la muestra distribuye Cauchy en lugar de normal están dadas por<sup>16</sup> por la figura 8.

El siguiente código se empleó para realizar las dos gráficas anteriores de donde simulo de una *distribución Cauchy*:

```

1
2 # Construyo la base multipanel para los diferentes a, asumiendo que  $Y_{\{i\}}^{\{0\}}$  sigue una
   distribuci n Cauchy
3
4 base_multipanel_cauchy = base_para_graficas(t2, n_vector2, distro = "cauchy", truth_delta =
   3)
5
6 # II. ----
7
8 base_filtrada_cauchy = base_multipanel_cauchy %>%
9   filter(tama o > 20)
10  # filter(tama o < 10000)
11
12 base_filtrada_cauchy2 = base_multipanel_cauchy %>%
13   filter(tama o > 10) %>%
14   filter(tama o < 2000)
15
16 # Gr ficas de las funciones de densidad para la simulaci n de los delta estimados
17 # y de los a, para el caso de un  $Y_{\{i\}}^{\{0\}} = \text{cauchy}(0, 1)$ 
18 grafica2_a = compracion_densidades_grafica(t = t2, df = base_filtrada_cauchy2, var = a);
   grafica2_a
19 grafica2_delta_est = compracion_densidades_grafica(t = t2, df = base_filtrada_cauchy, var =
   delta_est); grafica2_delta_est

```

Como se puede observar de 8 claramente, no se satisface el *teorema del limite central* cuando se tiene una muestra que *distribuye Cauchy*.

El teorema del limite central no se satisface principalmente porque la *distribución Cauchy* es una *distribución patológica*, en el sentido que no es una distribución que tenga las propiedades usuales que uno esperaría en una distribución. En específico, la distribución Cauchy, tiene media infinita y varianza

<sup>15</sup>A medida que el tamaño de las muestras aumente, es más probable que surgan valores atípicos muy extremos por lo que es más probable encontrar valores atípicos bastante extremos en los estimadores MCO para muestras de tamaño 1000 en comparación con las muestras de tamaño 20

<sup>16</sup>Por las instrucciones del taller, se emplearon muestras de tamaño 20, 100, 1000. Las muestras de tamaño  $n = 10$  no se emplearon por los problemas de multicolinealidad comentados anteriormente



Figura 8: Distribuciones asintóticas de  $a_{k,n}$  en donde claramente **no** se satisface el teorema de limite central para una muestra que distribuye Cauchy.

infinita, lo cuál hace que sea una distribución que tenga características muy particulares, entre ellas que pueda tener valores extremos muy atípicos y además que no satisfaga los teoremas clásicos de teoría asintótica como la ley débil de los grandes números y el teorema del límite central. Si se mira con detalle, el teorema del límite central está planteado para una distribución en la que se conozca la media poblacional y varianza poblacional de la variable aleatoria en la que se está aplicando el teorema. Como en una variable aleatoria, dicha media poblacional y varianza poblacional no existe, o más específicamente son infinitas, entonces no es posible que se satisfaga el teorema del límite central porque dos de las propiedades necesarias que debe tener una variable aleatoria para satisfacer el teorema de límite central, en particular tener valor esperado y varianza finita, no se satisface. En particular, cuando se tiene una variable aleatoria Cauchy, se sabe que en lugar de que su media muestral, debidamente normalizada, en lugar de converger a una distribución normal con varianza conocida, distribuye también como una variable aleatoria Cauchy<sup>17</sup>

Ahora bien, si empleáramos *inferencia convencional* para parámetros MCO a una muestra que distribuya Cauchy, se van a obtener conclusiones incorrectas a partir de dicha inferencia que se realice. Esto se debe a que la inferencia convencional que uno emplea cuándo hace estimaciones por MCO, solo es válida si se satisface la teoría asintótica del estimador MCO, y la distribución del estimador MCO es la esperada por lo que dictamina la teoría asintótica, en particular, dicha inferencia convencional que se usa cuando se emplea el estimador MCO depende crucialmente de que se satisfaga el teorema del límite central. Ahora bien, como en una muestra que distribuya Cauchy, en general no se cumple la teoría asintótica convencional, y en particular no se satisface el teorema del límite central, toda inferencia convencional del estimador MCO que se emplee bajo dicha muestra Cauchy será incorrecta dado que, supuestos importante como el cumplimiento del teorema del límite central no se da, y eso que la distribución asintótica del estimador de MCO deje de ser la esperada por teoría asintótica.

Finalmente, suponga que usted sabe que el programa implementado por el ministerio tiene el objetivo de nivelar los conocimientos de los estudiantes que en el participan. Por dicha razón, los resultados de las pruebas de los estudiantes que no participan en el programa son mas volátiles.

Una manera de modelar esta observación es a través del modelo:

$$Y_i = 3 * D_i + v_i$$

<sup>17</sup>Es decir, la media muestral de variables Cauchy i.i.d. también distribuye como una variable aleatoria Cauchy.

donde  $v_i \sim N(0, 2 - D_i)$ .

e) Para este nuevo modelo, realice el siguiente procedimiento:

- I. Simule 1000 muestras  $\{(Y_i(0), Y_i(1), D_i)\}_{i=1}^{1000}$ . Para cada muestra, recupere los intervalos de confianza clásicos al 90 %, 95 % y 99 % obtenidos para  $\delta$  y codifique en una matriz si dicho intervalo contiene o no a  $\delta$ .

• Solución:

Se tiene que el siguiente modelo representa las observaciones para el caso expuesto en este inciso:

$$Y_i = 3D_i + v_i \quad (8)$$

Ahora bien, como conocemos que  $v_i \sim \mathcal{N}(0, 2 - D_i)$ , entonces se puede decir que una representación de la distribución de  $Y_i$  estaría dada por:

$$Y_i = 3D_i + \mathcal{N}(0, 2 - D_i)$$

Por tanto, si  $D_i = 0$ , entonces  $Y_i = Y_i(0) = 0 + \mathcal{N}(0, 2)$ , lo que hace que:

$$Y_i(0) \sim \mathcal{N}(0, 2)$$

y si  $D_i = 1$ , entonces  $Y_i = Y_i(1) = 3 + \mathcal{N}(3, 1)$ , lo que hace que:

$$Y_i(1) \sim \mathcal{N}(3, 1)$$

Dicho lo anterior, se simulan 1000 muestras de tamaño  $n = 1000$  cada una, y se recuperan los intervalos de confianza clásicos al 90 %, 95 % y 99 % obtenidos para  $\delta$  y se codifican en una matriz que informa si dicho intervalo contiene o no  $\delta$ . El código para realizar ello es el siguiente:

```

1
2 # I. ----
3
4 # Tamaños de muestra
5 n20 = 20
6 n100 = 100
7 n200 = 200
8 n500 = 500
9 n1000 = 1000
10 n10000 = 10000
11
12 hetero20_no_robustos = delta(t = 1000, n = 20, distro = "normal", hetero = T, varianza_
   y0 = 2, varianza_y1 = 1, robustos = F)
13 hetero100_no_robustos = delta(t = 1000, n = 100, distro = "normal", hetero = T, varianza
   _y0 = 2, varianza_y1 = 1, robustos = F)
14 hetero200_no_robustos = delta(t = 1000, n = 200, distro = "normal", hetero = T, varianza
   _y0 = 2, varianza_y1 = 1, robustos = F)
15 hetero500_no_robustos = delta(t = 1000, n = 500, distro = "normal", hetero = T, varianza
   _y0 = 2, varianza_y1 = 1, robustos = F)
16 hetero1000_no_robustos = delta(t = 1000, n = 1000, distro = "normal", hetero = T,
   varianza_y0 = 2, varianza_y1 = 1, robustos = F)
17 hetero10000_no_robustos = delta(t = 1000, n = 10000, distro = "normal", hetero = T,
   varianza_y0 = 2, varianza_y1 = 1, robustos = F)
18
19 n_vect_hetero = c(n20, n100, n200, n500, n1000, n10000)
20
21 # Dataframe con el tamaño de muestra y la media y varianza del
22 # estimador de delta para los diferentes tamaños de muestra
23 media_varianza_hetero = media_varianza_delta(t = 1000, n_vect_hetero, distro = "normal",
   hetero = T, varianza_y0 = 2, varianza_y1 = 1, robustos = F); glimpse(media_varianza
   _hetero)
24
25 # Nota: Se observa que a pesar de la heterocedasticidad en  $Y_{\{i\}}$  hay insesgadez
26 # En la estimación de delta
27 # De igual forma, la varianza muestral también disminuye a medida que aumenta
28 # la muestra (independiente de que haya heterocedasticidad en  $Y_{\{i\}}$ )
29
30 # Construir una función que me compute los intervalos de confianza
31 # Intervalos de confianza: Ya sea para intervalos clásicos o para intervalos robustos
32 # a la heterocedasticidad como los calculados por la matriz de
   varianzas y covarianzas Huber-White
33
34 inter_confianza = function(df_estimaciones, delta_true, int_conf){
35   norm_inf = qnorm((1 - int_conf)/2)
36   norm_sup = qnorm((1 - int_conf)/2 + int_conf)
37   lim_inf = df_estimaciones$delta_est + norm_inf * sqrt(df_estimaciones$var_ols)
38   lim_sup = df_estimaciones$delta_est + norm_sup * sqrt(df_estimaciones$var_ols)
39   contiene_o_no = c()
40   for (i in 1:length(lim_inf)){
41     if ((lim_inf[i] < delta_true) && (delta_true < lim_sup[i])){
42       contiene_o_no = append(contiene_o_no, 1)

```

```

43     }else{
44         contiene_o_no = append(contiene_o_no, 0)
45     }
46 }
47 # Variable del dataframe
48 df = data.frame(lim_inf, lim_sup, contiene_o_no)
49 return(df)
50 }
51
52 # Intervalos de confianza del 90 %
53 int_conf1000_no_robustos_90 = inter_confianza(hetero1000_no_robustos, delta_true = 3,
54         int_conf = 0.9)
55
56 # Intervalos de confianza del 95 %
57 int_conf1000_no_robustos_95 = inter_confianza(hetero1000_no_robustos, delta_true = 3,
58         int_conf = 0.95)
59
60 # Intervalos de confianza del 99 %
61 int_conf1000_no_robustos_99 = inter_confianza(hetero1000_no_robustos, delta_true = 3,
62         int_conf = 0.99)

```

II. ¿Qué porcentaje de los intervalos de cada nivel contiene al parámetro verdadero? Presente sus resultados en una tabla ¿A qué se deben estos resultados? ¿Qué le sugiere esto sobre su forma de computar los intervalos?

- Solución:

Para encontrar el porcentaje de los intervalos de cada nivel de significancia que contiene al parámetro verdadero, se empleo el siguiente código:

```

1
2 # II. ----
3
4 # Porcentaje de los intervalos de confianza del 90 % que contienen el par metro
5 porcentaje_no_robustos_90 = sum(int_conf1000_no_robustos_90$contiene_o_no)/nrow(int_
6     conf1000_no_robustos_90) * 100; porcentaje_no_robustos_90
7
8 # Porcentaje de los intervalos de confianza del 95 % que contienen el par metro
9     verdadero
10 porcentaje_no_robustos_95 = sum(int_conf1000_no_robustos_95$contiene_o_no)/nrow(int_
11     conf1000_no_robustos_95) * 100; porcentaje_no_robustos_95
12
13 # Porcentaje de los intervalos de confianza del 99 % que contienen el par metro
14     verdadero
15 porcentaje_no_robustos_99 = sum(int_conf1000_no_robustos_99$contiene_o_no)/nrow(int_
16     conf1000_no_robustos_99) * 100; porcentaje_no_robustos_99

```

Los resultados se pueden observar mejor en la tabla 1<sup>18</sup>.

Cuadro 1: Tabla que muestra el porcentaje de los intervalos clásicos que contienen el parámetro poblacional  $\delta = 3$  a diferentes niveles de significancia

Nivel de significancia	Porcentaje de los intervalos de cada nivel que contiene al parámetro verdadero
90 %	93.9 %
95 %	98 %
99 %	99.8 %

Lo que se observa, es que en presencia de heterocedasticidad para los errores del modelo, los intervalos de confianza clásicos rechazan menos veces de las que deberían al nivel de significancia propuesta la hipótesis nula, es decir, son más amplios de lo que deberían ser, y por tanto, hay más muestras que contienen el parámetro verdadero  $\delta = 3$  de las que deberían dado el nivel de significancia escogido, y esto ocurre para todos los niveles de significancia. Pareciera que los intervalos de confianza son demasiado amplios de manera sistemática sin importar el nivel de significancia. Esto se debe principalmente, a que la matriz de varianzas y covarianzas, bajo el supuesto de homocedasticidad, no es la adecuada a la representar la matriz de varianzas y covarianzas asintótica del estimador MCO cuando hay heterocedasticidad en la muestra.

III. Repita los numerales I y II, esta vez empleando errores estándar de White (robustos) en la construcción de sus intervalos. Concluya.

- Solución:

<sup>18</sup>Es importante notar que se simuló en total  $t = 1000$  muestras, donde cada muestra tenía un tamaño de muestra de  $n = 1000$ . Lo que permite concluir que, los resultados obtenidos tanto en el inciso II. como en el inciso III. (que utilizó la misma metodología) son bastante robustos y confiables, y muestran claramente el comportamiento tanto de los intervalos clásicos como de los intervalos construidos con errores robustos ante la presencia de heterocedasticidad en la muestra.

Dado lo encontrado en los incisos I. y II. del punto e., se propone a repetir el mismo procedimiento que el anterior pero esta vez utilizando una *matriz de varianzas y covarianzas* correspondientes a *errores robustos a la heterocedasticidad tipo Huber-White*. Para ello, se emplea el siguiente fragmente de código:

```

1
2 # III. ----
3
4 hetero20_robustos = delta(t = 1000, n = 20, distro = "normal", hetero = T, varianza_y0 =
  2, varianza_y1 = 1, robustos = T)
5 hetero100_robustos = delta(t = 1000, n = 100, distro = "normal", hetero = T, varianza_y0
  = 2, varianza_y1 = 1, robustos = T)
6 hetero200_robustos = delta(t = 1000, n = 200, distro = "normal", hetero = T, varianza_y0
  = 2, varianza_y1 = 1, robustos = T)
7 hetero500_robustos = delta(t = 1000, n = 500, distro = "normal", hetero = T, varianza_y0
  = 2, varianza_y1 = 1, robustos = T)
8 hetero1000_robustos = delta(t = 1000, n = 1000, distro = "normal", hetero = T, varianza_
  y0 = 2, varianza_y1 = 1, robustos = T)
9 hetero10000_robustos = delta(t = 1000, n = 10000, distro = "normal", hetero = T,
  varianza_y0 = 2, varianza_y1 = 1, robustos = T)
10
11 n_vect_hetero = c(n20, n100, n200, n500, n1000, n10000)
12
13 # Dataframe con el tamaño de muestra y la media y varianza del
14 # estimador de delta para los diferentes tamaños de muestra
15 media_varianza_hetero = media_varianza_delta(t = 1000, n_vect_hetero, distro = "normal",
  hetero = T, varianza_y0 = 2, varianza_y1 = 1, robustos = T); glimpse(media_varianza
  _hetero)
16
17 # III. Cálculo de intervalos de confianza cuando se tienen errores robustos ----
18
19 # Intervalos de confianza del 90 %
20 int_conf1000_robustos_90 = inter_confianza(hetero1000_robustos, delta_true = 3, int_conf
  = 0.9)
21
22 # Intervalos de confianza del 95 %
23 int_conf1000_robustos_95 = inter_confianza(hetero1000_robustos, delta_true = 3, int_conf
  = 0.95)
24
25 # Intervalos de confianza del 99 %
26 int_conf1000_robustos_99 = inter_confianza(hetero1000_robustos, delta_true = 3, int_conf
  = 0.99)
27
28 # III. Porcentaje de los intervalos de confianza cuando se tienen errores robustos ----
29
30 # Porcentaje de los intervalos de confianza del 90 % que contienen el parámetro
  verdadero
31 porcentaje_robustos_90 = sum(int_conf1000_robustos_90$contiene_o_no)/nrow(int_conf1000_
  robustos_90) * 100; porcentaje_robustos_90
32
33 # Porcentaje de los intervalos de confianza del 95 % que contienen el parámetro
  verdadero
34 porcentaje_robustos_95 = sum(int_conf1000_robustos_95$contiene_o_no)/nrow(int_conf1000_
  robustos_95) * 100; porcentaje_robustos_95
35
36 # Porcentaje de los intervalos de confianza del 99 % que contienen el parámetro
  verdadero
37 porcentaje_robustos_99 = sum(int_conf1000_robustos_99$contiene_o_no)/nrow(int_conf1000_
  robustos_99) * 100; porcentaje_robustos_99
38
39 # Conclusión: Los errores robustos calculados con matriz de varianzas y covarianzas
  Huber-White,
40 #
41 # sirven cuando hay heterocedasticidad en los errores de modelo
42 # Por el contrario, si se utilizan intervalos de confianza clásicos
43 # sin corregir por heterocedasticidad, se observa que los intervalos
44 # de confianza clásicos de manera sistemática para los niveles de
45 # significancia sobreestiman, es decir calculan intervalos de
46 # confianza más amplios de los que en realidad deberían computarse
  dada la presencia de la heterocedasticidad

```

En dicho código, se generan los intervalos de confianza pero ahora contruidos con una matriz de varianzas y covarianzas correspondiente a *errores robustos a la heterocedasticidad tipo Huber-White* y se calcula cuántos de éstos intervalos contienen el parámetro poblacional  $\delta = 3$

Los resultados se pueden observar mejor en la tabla 2.

Cuadro 2: Tabla que muestra el porcentaje de los intervalos contruidos con errores robustos que contienen el parámetro poblacional  $\delta = 3$  a diferentes niveles de significancia

Nivel de significancia	Porcentaje de los intervalos de cada nivel que contiene al parámetro verdadero
90 %	89.7 %
95 %	94.9 %
99 %	99.6 %

Lo que se puede concluir, es que al emplear intervalos construidos con errores robustos basados en una matriz de varianzas y covarianzas *Huber-White*, es posible observar que el porcentaje de intervalos que contienen el parámetro poblacional  $\delta = 3$  es muy cercano al nivel de significancia propuesto. Lo que quiere decir, es que al emplear una matriz de varianzas y covarianzas *Huber-White*, claramente los intervalos de confianza tienen una longitud que corresponde mucho mejor al nivel de significancia propuesto por lo que claramente éstos intervalos de confianza son mejores que los intervalos de confianza clásicos como se puede evidenciar fácilmente de las tablas 1 y 2. Lo anterior, se explica por el hecho de que la matriz de varianzas-covarianzas *Huber-White* está diseñada para capturar la posible heterocedasticidad que puede presentar la muestra por lo que se puede ver como una generalización de la matriz de varianzas y covarianzas homocedásticas, y por tanto, en presencia de heterocedasticidad, arroja mejores intervalos de confianza.

## Punto doctorado

Usted se encuentra cursando el curso de seminario de investigación doctoral. En las primeras semanas le piden que presente algunas relaciones de causalidad que tenga pensado abordar en su disertación. Su profesora de seminario lo invita a que avance en el planteamiento de sus ideas a través del uso de gráficos acíclicos dirigidos (Directed Acyclic Graph - DAG).<sup>19</sup> En breve, los DAGs se componen de dos elementos: flechas y variables. Las flechas entre variables indican causalidad en el sentido en el que esta apuntando la flecha. Un **Camino** es cualquier conexión entre dos variables realizada por flechas, sin importar su dirección o si existen otras variables intermedias. Un **Camino por la puerta trasera** de la variable  $X$  a la variable  $Y$  es un camino que empieza con una flecha dirigida hacia  $X$ . Un **Colisionador** es una variable a la que apuntan dos flechas en un camino.

- a) Usando puntos negros para graficar las variables observadas y circunferencias para graficar las variables no observadas ( $X$ ,  $Z$  y  $Y$ ), realice los DAGs que representan:
  - i)  $X$  tiene una relación de causalidad hacia  $Y$ . A su vez,  $Z$  tiene una relación de causalidad hacia  $X$  y otra hacia  $Y$ .
  - ii)  $X$  tiene una relación de causalidad hacia  $Z$ . A su vez,  $Z$  tiene una relación de causalidad hacia  $Y$ .
  - iii) La variable no observada:  $U$  tiene una relación de causalidad hacia  $Z$  y otra hacia  $Y$ . A su vez,  $X$  tiene una relación de causalidad hacia  $Y$  y otra hacia  $Z$ .
  - iv) La variable  $X$  tiene una relación de causalidad hacia  $Z$  y otra hacia  $Y$ . A su vez,  $Y$  tiene una relación de causalidad hacia  $Z$ .
  - v) La variable no observada:  $U$  tiene una relación de causalidad hacia  $Z$  y otra hacia  $Y$ . A su vez,  $X$  tiene una relación de causalidad hacia  $Y$ , mientras que  $Z$  tiene una relación de causalidad hacia  $X$ .
- b) Ahora usted quiere implementar un modelo de regresión lineal simple para evaluar los efectos causales de  $X$  sobre  $Y$ . Con este objetivo, usted está determinando la conveniencia de incluir la variable  $Z$  como variable de control en su modelo. Argumente, para cada caso representado por los DAGs, si la inclusión de  $Z$ , como control en el modelo, es apropiada. Para sus respuestas, tenga en cuenta que un variable es un mal control si: bloquea caminos causales entre  $X$  y  $Y$  o abre otros caminos que no son causales entre  $X$  y  $Y$ .
 

*(Pista: Bloquear un camino es equivalente a controlar por variables que no son Colisionadoras o no controlar por las Colisionadoras. Abrir un camino es equivalente a controlar por variables que son Colisionadoras o no controlar por no Colisionadoras)*
- c) Finalmente, usted sabe que la variable  $X$  es independiente de los resultados potenciales de  $Y$  y de  $Z$ . Por otro lado, usted sospecha que  $Z$  podría ser un mal control y que este hecho puede generarse porque la variable  $Z$  es, a su vez, una variable de resultado en su modelo. Plantee una forma de aproximarse empíricamente para evaluar si este hecho puede estar ocurriendo

<sup>19</sup>Si no se encuentra familiarizado con los DAGs, una excelente referencia y explicación se encuentra en el tercer capítulo del libro de Scott Cunningham: Causal Inference: The mixtape, del año 2021. <https://mixtape.scunning.com/>



## Anexo: Código segundo ejercicio

A continuación se anexa todo el código empleado para resolver el *segundo ejercicio* del taller:

```
1 // Taller 1
2 use "/Users/federicoduenas/Desktop/Econometri a_ Avz_/taller 1/manufacturaCol.dta"
3
4 foreach var of varlist producto_acum producto costos {
5     gen ln_`var' = ln(`var'+1)
6 }
7
8
9
10 // relacion de log_c con log_n controlando por log_y
11 // usando Partialling out:
12
13 reg ln_producto_acum ln_producto
14
15 gen aux_1 = 9.256792 + 0.1120002*ln_producto
16
17 // para tener la parte de N que no es explicada por Y
18 gen e_1 = ln_producto_acum - aux_1
19
20 // ahora para tener los C_i no explicada por Y_i
21
22 reg ln_costos ln_producto
23
24 gen aux_2_1 = -.019997 + .0052894*ln_producto
25
26 // los residuales: parte de C_i no explicada por Y_i
27 gen e_1_2 = ln_costos - aux_2_1
28
29 // sactter de la relaci n C_i con N_i una vez se removi el efecto de Y_i
30 ***#
31
32 scatter e_1_2 e_1 || lfit e_1_2 e_1
33
34 /*
35
36
37 // ahora si calculo la relacion entre N y C, controlando por Y
38
39 reg ln_costos e_1
40
41 // Que, por teorema de Waugh-Frisch-Lovell, es equivalente al par metro de log_N de la reg lineal:
42 reg ln_costos ln_producto_acum ln_producto
43
44 // gr fica de la relaci n de C y N:
45 scatter ln_costos e_1 || lfit ln_costos e_1
46 */
47 // ahora la relacion de log_C y log_Y controlando por log_N
48
49 reg ln_producto ln_producto_acum
50
51 gen aux_2 = .3859203 + .4743226*ln_producto_acum
52
53 gen e_2 = ln_producto - aux_2
54
55
56 // ahora remuevo el efecto de N_i sobre C_i
57
58 reg ln_costos ln_producto_acum
59 gen aux_2_2 = -.013807 + .0020865*ln_producto_acum
60 gen e_2_2 = ln_costos - aux_2_2
61
62
63 // sactter de la relaci n de C y Y removiendo el efecto de N_i
64
65 scatter e_2_2 e_2 || lfit e_2_2 e_2
66
67 /*
68 reg ln_costos e_2
69
70 reg ln_costos ln_producto_acum ln_producto
71
72
73 scatter ln_costos e_2 || lfit ln_costos e_2
74
75 // estimar la ecuaci n del modelo:
76
77 reg ln_costos ln_producto_acum ln_producto
78 */
79 reg ln_costos ln_producto_acum ln_producto
80
81
82 matlist e(b)
83 matlist e(V)
```



## Anexo: Código tercer ejercicio

A continuación se anexa todo el código empleado para resolver el *tercer ejercicio* del taller:

```
1 library(tidyverse)
2
3 # Definir la semilla para reproducibilidad en los resultados
4 set.seed(12345)
5
6 # Punto 3 ----
7
8 # Ejercicio a. ----
9
10 # Switching regression:
11 ##  $y_{\{i\}} = y_{\{i\}}^{\{1\}} * D_{\{i\}} + (1 - D_{\{i\}}) * y_{\{i\}}^{\{0\}}$ 
12 ##  $y_{\{i\}} = y_{\{i\}}^{\{0\}} + (y_{\{i\}}^{\{1\}} - y_{\{i\}}^{\{0\}}) * D_{\{i\}}$ 
13 ##  $y_{\{i\}} = \text{delta} * D_{\{i\}} + \text{epsilon}_{\{i\}}$ 
14
15 # Donde el par metro causal:
16 ##  $\text{delta} = E[y_{\{i\}}^{\{1\}} | D_{\{i\}} = 1] - E[y_{\{i\}}^{\{0\}} | D_{\{i\}} = 0]$ 
17 ##  $\text{delta} = 3 - 0$ 
18 ##  $\text{delta} = 3$ 
19
20 # Ejercicio b. ----
21
22 # Nota: Dado las condiciones para simular  $y_0 = y_{\{i\}}^{\{0\}}$  y  $y_1 = y_{\{i\}}^{\{2\}}$ ,
23 # uno esperar a que el SDO = 3 y por tanto el delta = 3
24 # Ademas, dado que el SB = 0 (teniendo en cuenta que  $y_0$  y  $y_1$  son independientes de D)
25 # se tiene que delta representa el impcto causal del programa del ministerio
26
27 # I. ----
28
29 # Definición de par metros para la simulación
30
31 ## Número de muestras:
32 t = 100
33
34 ## Tamaño de muestras:
35
36 n10 = 10
37 n20 = 20
38 n50 = 50
39 n100 = 100
40 n500 = 500
41 n1000 = 1000
42
43 #-----
44 # Funciones auxiliares:
45 #-----
46
47 # Cálculo del SDO: (SDO := Simple difference of outcomes) (Función auxiliar)
48 SDO_calculo = function(Y, D){
49   # Por construcción
50   ##  $Y = Y_0$  Si  $D = 0$ 
51   ##  $Y = Y_1$  Si  $D = 1$ 
52   # Nota: Y y D son vectores del mismo tamaño
53   # Inicialización de contadores
54   count0 = 0
55   count1 = 0
56   # Inicialización de la suma para cálculo del SDO
57   suma0 = 0
58   suma1 = 0
59   # Iteración a través de la muestra
60   for (i in 1:length(Y)){
61     if (D[i] == 0){
62       suma0 = suma0 + Y[i]
63       count0 = count0 + 1
64     }else{
65       suma1 = suma1 + Y[i]
66       count1 = count1 + 1
67     }
68   }
69   SDO = (suma1/count1) -(suma0/count0) # SDO: Simple diferencia de muestras
70   return(SDO)
71 }
72
73 # Construcción matriz X para un modelo de regresión lineal con constante: (Función auxiliar)
74 ols_X = function(...){
75   # ... contiene los regresores necesarios para construir la matriz X
76   regresores = list(...)
77   # Nota: Cada regresor es un vector con el mismo número de observaciones n
78   n = length(regresores[[1]]) # No importa que se tome el primer regresor dado que todos los
79   # regresores tienen el mismo tamaño
80   const = rep(1, times = n)
81   X = cbind(const, ...)
82   return(X)
83 }
84
85 # Construcción del estimador  $\hat{\sigma}^2$ : (Función auxiliar)
86 estimador_sigma = function(y, X, beta_est){
87   # Variables:
```

```

87  ## y: es el vector que contiene la variable dependiente
88  ## X: es la matriz que contiene las variables regresoras
89  ## beta_est: es el vector de los par metros estimados por OLS
90  e = y - X %*% beta_est # e es el vector de residuales
91  k = ncol(X) # El n mero de par metros a estimar
92  n = length(y) # El n mero de observaciones en la muestra
93  sigma_est = (t(e) %*% e) / (n - k) # estimador de sigma (varianza del t rmino de error)
94  # El resultado de la funci n es: sigma_est (que es un escalar)
95  return(sigma_est)
96 }
97
98 # Construcci n del estimador para la matriz \Hat{Q}_{XX}^{-1}: (Funci n auxiliar)
99 estimador_Q_xx_inv = function(X){
100   n = nrow(X) # donde n es el tama o de la muestra
101   Q = solve((1/n) * (t(X) %*% X)) # Q es el estimador que se est buscando
102   # Q es una matriz de tama o k x k, donde k es en n mero de par metros a estimar
103   # La funci n me retorna a una matriz Q de tama o k x k
104   return(Q)
105 }
106
107 # Construcci n del estimador de la matriz de varianzas y covarianzas
108 # Huber-White: (Funci n auxiliar)
109 Huber_White = function(y, X, beta_est){
110   # Variables:
111   ## y: es el vector que contiene la variable dependiente
112   ## X: es la matriz que contiene las variables regresoras
113   ## beta_est: es el vector de los par metros estimados por OLS
114   e = y - X %*% beta_est # e es el vector de residuales
115   k = ncol(X) # El n mero de par metros a estimar
116   n = length(y) # El n mero de observaciones en la muestra
117   # Construcci n de (X^{'}X)^{-1}
118   mat_X = solve(t(X) %*% X)
119   # Matriz intermedia de los errores robustos White (que se encuentra por medio de una suma)
120   # Inicializo la matriz como una matriz de ceros
121   mat_intermedia = matrix(0, nrow = k, ncol = k)
122   # El for se dise a para llenar la matriz mat_intermedia
123   for (i in 1:n){
124     mat_intermedia = mat_intermedia + (e[i])^2 * (X[i, ] %*% t(X[i, ]))
125   }
126   # mat_final me dal estimador de la matriz de varianzas y covarianzas de Huber-White
127   mat_final = mat_X %*% mat_intermedia %*% mat_X
128   return(mat_final)
129 }
130
131 # y = rnorm(1000)
132 # D = rnorm(1000, mean = 2)
133 # beta = c(1, 2)
134 # X = ols_X(D)
135 #
136 # Huber_White(y, X, beta)
137
138
139 # l estimador para la matriz \Hat{Q}_{XX}^{-1}: (Funci n auxiliar)
140
141 ##-----
142 # delta: Es la funci n principal (m s importante) de todo el c digo.
143 # Con base en la funci n delta es que se deduce todo lo dem s
144 # en el c digo.
145 # Es una funci n bastante general que contempla todas las situaciones
146 # que puedan surgir en el c digo. En ese caso, contempla la simulaci n
147 # de la Cauchy y de una muestra heteroced stico donde la volatilidad
148 # de los tratados es diferente a la volatilidad de los no tratados
149 # Nota: Revisar toda el enunciado del ejercicio primero para entender mejor
150 # la l gica de la funci n
151 ##-----
152
153 # Funci n para estimar el efecto tratamiento (delta) en una base de datos
154 delta = function(t, n, distro, hetero = F, varianza_y0 = 1, varianza_y1 = 1, robustos = F){
155   # Defino la semilla para reproducibilidad del resultado
156   set.seed(5678)
157   # Definici n de variables:
158   ## t: n mero de muestras a simular
159   ## n: tama o de las muestras
160   ## distro: tipo de distribuci n con el que se va a simular Y_{i}^{0}
161   ## Nota: Si distro == "normal", entonces los par metros opcionales
162   ## empiezan a tomar relevancia
163   ## hetero: Para saber si la volatilidad entre tratados y no tratados es diferente
164   ## varianza_y0: El valor de la varianza de la distribuci n normal con la que se simula y_{i}^{0}
165   ## varianza_y1: El valor de la varianza de la distribuci n normal con la que se simula y_{i}^{1}
166   ## robustos: Solo se activa si se escoge la opci n hetero == T y adem s robustos == T
167   ## Permite estimar errores robustos usando la matriz de varianzas y covarianzas Huber-
168   White
169   # Consideraciones adicionales:
170   # Condici n l gica para saber si y_0 sigue una distribuci n normal est ndar o una
171   # distribuci n Cauchy est ndar
172   if (distro == "normal"){
173     # Hay dos posibles casos en caso de que la distribuci n sea normal o no:
174     # 1. Homoced sticidad entre tratados y no tratados
175     ## En caso que la volatilidad de y_{i}^{0} y y_{i}^{1} sea la misma hetero == F
176     # 2. Heterocedastidad entre tratados y no tratados (diferente volatilidad dependiendo si fueron
177     # tratados o no)
178     ## En caso que la volatilidad de y_{i}^{0} y y_{i}^{1} sea diferente hetero == T

```

```

176 if (hetero == F){
177   # df: Dataframe que almacena el SDO y el delta_estimado por OLS
178   # El par metro delta es el para metro asociado a la asignaci n a tratamiento (D)
179   df = data.frame(SDO = double(), delta_est = double(), sigma = double(), var_tlc = double(),
180                   var_ols = double())
181   # Meter un for para generar las t diferentes muestras
182   for (muestra in 1:t){
183     # Simulaci n de los outcomes potenciales y de la variable trataminto
184     y_0 = rnorm(n, mean = 0, sd = 1) # Simulaci n del outcome potencial de ausencia de
185     tratamiento
186     y_1 = y_0 + 3 # Simulaci n del outcome potencial de presencia de tratamiento
187     D = rbinom(n, 1, prob = 0.3) # Simulaci n de una variable Bernoulli con una probabilidad
188     de xito de 0.3
189     # Modelo causal de Rubin
190     y = y_0 + (y_1 - y_0) * D # y es un vector Nx1, donde n es el n mero de observaciones (
191     tama o de muestra)
192     # Corregir (Debo calcular es el SDO)
193     SDO = SDO_calculo(y, D) # C lculo del SDO utilizando la funci n SDO_calculo
194     # Estimaci n del par metro delta mediante una regresi n lineal con constante
195     # X es una matriz NxK, donde n es el n mero de observaciones y K el n mero de
196     para etros
197     X = ols_X(D) # D es la variable tratamiento
198     beta_ols = solve(t(X) %*% X) %*% t(X) %*% y # beta_ols es un vector de Kx1, donde K es el
199     n mero de par metros
200     # Ahora, se calcularon 3 par metros adicionales:
201     # el estimador de sigma, el estimador de var(sqrt(n) (delta_{OLS} - delta) y el estimador
202     de var(delta_{OLS}))
203     ## estimador de sigma:
204     sigma = estimador_sigma(y, X, beta_ols)
205     ## estimador de \Hat{Q}_{XX}^{-1}:
206     Q_mat = estimador_Q_xx_inv(X)
207     ## estimador de var_tlc = var(sqrt(n) (delta_{OLS} - delta))
208     var_tlc = sigma * Q_mat[[2,2]]
209     ## estimador de var_ols = var(delta_{OLS})
210     n = nrow(X)
211     var_ols = 1/n * sigma * Q_mat[[2,2]]
212     # Dataframe que contiene todos los par metros de inter s
213     df[muestra, ] = c(SDO, beta_ols[[2,1]], sigma, var_tlc, var_ols) # Extraigo el segundo
214     par metro que es el par metro delta
215   }
216 }else{
217   # df: Dataframe que almacena el SDO y el delta_estimado por OLS
218   # El par metro delta es el para metro asociado a la asignaci n a tratamiento (D)
219   df = data.frame(SDO = double(), delta_est = double(), sigma = double(), var_ols = double())
220   # Meter un for para generar las t diferentes muestras
221   for (muestra in 1:t){
222     # Simulaci n de los outcomes potenciales y de la variable trataminto
223     D = rbinom(n, 1, prob = 0.3) # Simulaci n de una variable Bernoulli con una probabilidad
224     de xito de 0.3
225     y_0 = rnorm(n, mean = 0, sd = sqrt(varianza_y0)) # Simulaci n del outcome potencial de
226     ausencia de tratamiento
227     y_1 = rnorm(n, mean = 3, sd = sqrt(varianza_y1)) # Simulaci n del outcome potencial de
228     presencia de tratamiento
229     # Modelo causal de Rubin
230     y = y_0 + (y_1 - y_0) * D # y es un vector Nx1, donde n es el n mero de observaciones (
231     tama o de muestra)
232     # Corregir (Debo calcular es el SDO)
233     SDO = SDO_calculo(y, D) # C lculo del SDO utilizando la funci n SDO_calculo
234     # Estimaci n del par metro delta mediante una regresi n lineal con constante
235     # X es una matriz NxK, donde n es el n mero de observaciones y K el n mero de
236     para etros
237     X = ols_X(D) # D es la variable tratamiento
238     beta_ols = solve(t(X) %*% X) %*% t(X) %*% y # beta_ols es un vector de Kx1, donde K es el
239     n mero de par metros
240     if (robustos == F){
241       # Ahora, se calcularon 3 par metros adicionales:
242       # el estimador de sigma, el estimador de var(sqrt(n) (delta_{OLS} - delta) y el estimador
243       de var(delta_{OLS}))
244       ## estimador de sigma:
245       sigma = estimador_sigma(y, X, beta_ols)
246       ## estimador de \Hat{Q}_{XX}^{-1}:
247       Q_mat = estimador_Q_xx_inv(X)
248       ## estimador de var_ols = var(delta_{OLS})
249       n = nrow(X)
250       var_ols = 1/n * sigma * Q_mat[[2,2]]
251       # Dataframe que contiene todos los par metros de inter s
252       df[muestra, ] = c(SDO, beta_ols[[2,1]], sigma, var_ols) # Extraigo el segundo
253       par metro que es el par metro delta
254     }else{
255       # Sigma, es provisional, luego se retira porque no tiene sentido para un estimador con
256       matriz de varianzas y covarianzas
257       # con errores robustos Huber-White (Dado la heterocedasticidad de los errores)
258       sigma = 0
259       # C lculo errores robustos Huber-White (me genera una matriz k x k, donde k es el
260       n mero de par metros)
261       var_ols = Huber_White(y, X, beta_ols)[[2,2]] # Extraigo la componente 2 de la matriz de
262       varianzas y covarianzas de Huber-White
263       df[muestra, ] = c(SDO, beta_ols[[2,1]], sigma, var_ols) # Extraigo el segundo
264       par metro que es el par metro delta
265     }
266   }
267 }
268 # Condicional final para eliminar la columna sigma si robustos == T

```

```

248     if (robustos == T){
249         df = df %>%
250             select(SD0, delta_est, var_ols)
251     }
252 }
253 }else if (distro == "cauchy"){
254     # df: Dataframe que almacena el SD0 y el delta_estimado por OLS
255     # El par metro delta es el par metro asociado a la asignación a tratamiento (D)
256     df = data.frame(SD0 = double(), delta_est = double())
257     for (muestra in 1:t){
258         # Simulación de los outcomes potenciales y de la variable tratamiento
259         y_0 = rcauchy(n, location = 0, scale = 1) # Simulación del outcome potencial de ausencia de
                tratamiento
260         y_1 = y_0 + 3 # Simulación del outcome potencial de presencia de tratamiento
261         D = rbinom(n, 1, prob = 0.3) # Simulación de una variable Bernoulli con una probabilidad de
                éxito de 0.3
262         # Modelo causal de Rubin
263         y = y_0 + (y_1 - y_0) * D # y es un vector Nx1, donde n es el número de observaciones (
                tamaño de muestra)
264         # Corregir (Debo calcular es el SD0)
265         SD0 = SD0_calculo(y, D) # Cálculo del SD0 utilizando la función SD0_calculo
266         # Estimación del par metro delta mediante una regresión lineal con constante
267         # X es una matriz NxK, donde n es el número de observaciones y K el número de par metros
268         X = ols_X(D) # D es la variable tratamiento
269         beta_ols = solve(t(X) %*% X) %*% t(X) %*% y # beta_ols es un vector de Kx1, donde K es el
                número de par metros
270         df[muestra, ] = c(SD0, beta_ols[[2,1]]) # Extraigo el segundo par metro que es el
                par metro delta
271     }
272 }
273 # La función delta retorna un df con el cálculo de la SD0
274 # y el par metro estimado delta, que proviene de una switching regression
275 return(df)
276 }
277
278 # Simulación
279 # delta(t, n10, distro = "normal") # Existe el riesgo de que haya
280 # multicolinealidad perfecta cuando se usa n10 = 10
281 # Puede generar un vector D = rep(0, times = 10)
282
283 delta20 = delta(t, n20, distro = "normal")
284 delta50 = delta(t, n50, distro = "normal")
285 delta100 = delta(t, n100, distro = "normal")
286 delta500 = delta(t, n500, distro = "normal")
287 delta1000 = delta(t, n1000, distro = "normal")
288
289 # II. ----
290
291 # n_vector es un vector que almacena los diferentes tamaños de muestra
292 n_vector = seq(from = 20, to = 1000, by = 10)
293
294 # Función que genera un dataframe con el tamaño de muestra,
295 # la media y la varianza del estimador de delta para cada
296 # tamaño diferente de muestra
297
298 media_varianza_delta = function(t, n_vector, distro, hetero = F, varianza_y0 = 1, varianza_y1 = 1,
                robustos = F){
299     # Variables:
300     ## t: Número de muestras que se van a generar por cada tamaño de muestra
301     ## n_vector: Variable que almacena los diferentes tamaños de muestras
302     ## distro, hetero, varianza_y0 y varianza_y1 son par metros definidos para la función delta
303     # df: Dataframe que almacena el tamaño de muestra,
304     # la media y la varianza del estimador de delta para cada tamaño de muestra
305     df = data.frame(tamaño = double(), media = double(), varianza = double())
306     # Llamo a la función delta para cada tamaño diferente de muestra
307     for (i in 1:length(n_vector)){
308         n_muestra = n_vector[i]
309         delta_n = delta(t, n_muestra, distro, hetero, varianza_y0, varianza_y1, robustos)$delta_est #
                delta_n es el vector de deltas por cada tamaño de muestra
310         df[i,] = c(n_muestra, mean(delta_n), var(delta_n))
311     }
312     return(df)
313 }
314
315 # Dataframe con el tamaño de muestra y la media y varianza del
316 # estimador de delta para los diferentes tamaños de muestra
317 media_varianza = media_varianza_delta(100, n_vector, distro = "normal"); glimpse(media_varianza)
318
319 # III. ----
320
321 grafica_propiedades = function(df, variable_y, y_intercepto = 0){
322     variable_y = ensym(variable_y)
323     graph = df %>%
324         ggplot(aes(x = tamaño, y = !!variable_y)) +
325         geom_line(color = "green", size = 1) +
326         geom_hline(yintercept = y_intercepto, color = "red") +
327         theme_light()
328     return(graph)
329 }
330
331 grafica_propiedades(media_varianza, media, y_intercepto = 3)
332 grafica_propiedades(media_varianza, varianza)

```

```

333
334 # Ejercicio c. ----
335
336 # I. ----
337
338 # Definici n de par metros para la simulaci n
339
340 ## N mero de muestras:
341 t2 = 1000
342
343 ## Tama o de muestras:
344
345 # No se puede trabajar con n10 = 10 porque no satisface la condici n de rango
346 # n10 = 10 Genera problemas de singularidad cuando se generan muchas muestras
347 # Porque por chance, se genera una muestra donde D = rep(0, times = 10)
348 n20 = 20
349 n100 = 100
350 n1000 = 1000
351 n5000 = 5000
352 n10000 = 10000
353 n20000 = 20000
354 n50000 = 50000
355 n100000 = 100000
356 n_vector2 = c(n20, n100, n1000, n5000, n10000, n20000, n50000, n100000)
357
358 # Funci n que calcula a = sqrt(n) (delta_est - truth_delta) para cada tama o de muestra n
359 calculo_a = function(t, n, distro, truth_delta){
360   # Variables:
361   ## t: n mero de muestras a simular
362   ## n: tama o de la muestra
363   ## truth_delta: verdadero valor del par metro delta (valor poblacional del par metro)
364   # df_estimados es el dataframe que contiene el delta estimado por SDO
365   # y por medio de regresi n
366   df_estimados = delta(t, n, distro)
367   # df_a es el dataframe que contiene
368   df_a = df_estimados %>%
369     mutate(a = sqrt(n) * (delta_est - truth_delta), tama o = rep(n, times = t))
370   return(df_a)
371 }
372
373 # Defino una funci n para crear un dataframe que se va a utilizar
374 # Para construir la gr fica multipanel.
375 base_para_graficas = function(t, n_vector, distro, truth_delta){
376   for (i in 1:length(n_vector)){
377     # Si i == 1 significa que estamos en la primera iteraci n
378     if (i == 1){
379       df_a_total = calculo_a(t, n_vector[i], distro, truth_delta = 3)
380     }else{
381       df_prov = calculo_a(t, n_vector[i], distro, truth_delta = 3)
382       df_a_total = bind_rows(df_a_total, df_prov)
383     }
384   }
385   # Retorna un dataframe con los diferentes valores de a para cada tama o de muestra listo
386   # para construir la gr fica multipanel
387   return(df_a_total)
388 }
389
390 # La base multipanel me tiene lista los diferentes a
391 # para cada tama o de muestra
392 base_multipanel = base_para_graficas(t2, n_vector2, distro = "normal", truth_delta = 3)
393
394 # II. ----
395
396 # Defino una funci n para crear la gr fica multipanel con el histograma y la funci n de densidad
397 # para cada tama o de muestra n
398 histogram_grid = function(df, titulo, x_lab, y_lab, num_bins = 30, y_upper_limit){
399   histog_grid = df %>%
400     ggplot(aes(a)) +
401     scale_y_continuous(limits = c(0, y_upper_limit)) +
402     geom_histogram(aes(y = ..density..), color = "black", bins = num_bins) +
403     geom_density(color = "green") +
404     facet_grid(cols = vars(tama o)) +
405     ggtitle(titulo) +
406     xlab(x_lab) +
407     ylab(y_lab) +
408     theme_light()
409 }
410
411 grafica_multi = histogram_grid(df = base_multipanel, titulo = "Gr fica multipanel para a con
diferentes tama os de muestra", x_lab = "Tama o de muestras (n)", y_lab = "", num_bins = 30,
y_upper_limit = 0.2); grafica_multi
412
413 # Gr fica que compara directamente las densidades
414 compracion_densidades_grafica = function(t, df, var){
415   # Variables:
416   ## t: n mero de muestras simuladas por tama o de muestra
417   ## df: df multipanel que tiene las simulaciones de las diferentes muestra por tama o de muestra
418   ## var: variable simulada de la cual se va a generar las funciones de densidad
419   # Nota:
420   # Modifico el dataframe multipanel para seleccionar solo
421   # las variables a y tama o que permiten hacer la gr fica de densidades
422   # df_mod es un dataframe modificado del dataframe multipanel

```

```

423 df_mod = df %>%
424   rename(valores_sim = {{ var }}, distribuciones = tama_o) %>%
425   mutate(distribuciones = as.character(distribuciones)) %>%
426   select(valores_sim, distribuciones)
427 # Simulo de una distribuci n normal est ndar porque quiero probar
428 # si se cumple o no el teorema del l mite central
429 df_normal = data.frame(valores_sim = rnorm(t), distribuciones = rep("Normal est ndar", times = t
  ))
430 # base_grafica ya es la base de datos lista para realizar las gr ficas de
431 # las densidades de las a y de una normal est ndar
432 base_grafica = df_mod
433 # base_grafica = bind_rows(df_mod, df_normal)
434 # Gr fica de densidades
435 density_comparacion = base_grafica %>%
436   ggplot(aes(x = valores_sim, color = distribuciones)) +
437   geom_density() +
438   theme_light() +
439   ggtitle("Comparaci n de funciones de densidad de a\n para muestras de diferentes tama os \n(
    todas simuladas)") +
440   ylab("Densidades") +
441   xlab("Distribuciones de a para diferente tama o de muestra")
442 return(density_comparacion)
443 }
444
445 # Visualizaci n de la base de datos con las variables simuladas
446 # delta_est y a
447 glimpse(base_multipanel)
448
449 base_filtrada = base_multipanel %>%
450   filter(100 < tama_o)
451
452 base_filtrada_a = base_multipanel %>%
453   filter(10 < tama_o) %>%
454   filter(tama_o < 2000)
455
456 # Gr ficas de las funciones de densidad para la simulaci n de los delta estimados
457 # y de los a, para el caso de un  $Y_{\{i\}}^{\{0\}} = \text{norm}(0, 1)$ 
458 grafica_a = compracion_densidades_grafica(t = t2, df = base_filtrada_a, var = a); grafica_a
459 grafica_delta_est = compracion_densidades_grafica(t = t2, df = base_filtrada, var = delta_est);
  grafica_delta_est
460
461 # Nota: Crear la otra gr fica mejor (script de simulaci n de df)
462
463 # Ejercicio d. ----
464
465 # Va a replicarse el ejercicio c. pero asumiendo que:
466 #  $Y_{\{i\}}^{\{0\}} \sim \text{cauchy}(0, 1)$  (Distribuci n Cauchy est ndar)
467
468 # I. ----
469
470 # Construyo la base multipanel para los diferentes a, asumiendo que  $Y_{\{i\}}^{\{0\}}$  sigue una
  distribuci n Cauchy
471
472 base_multipanel_cauchy = base_para_graficas(t2, n_vector2, distro = "cauchy", truth_delta = 3)
473
474 # II. ----
475
476 base_filtrada_cauchy = base_multipanel_cauchy %>%
477   filter(tama_o > 20)
478   # filter(tama_o < 10000)
479
480 base_filtrada_cauchy2 = base_multipanel_cauchy %>%
481   filter(tama_o > 10) %>%
482   filter(tama_o < 2000)
483
484 # Gr ficas de las funciones de densidad para la simulaci n de los delta estimados
485 # y de los a, para el caso de un  $Y_{\{i\}}^{\{0\}} = \text{cauchy}(0, 1)$ 
486 grafica2_a = compracion_densidades_grafica(t = t2, df = base_filtrada_cauchy2, var = a); grafica2_a
487 grafica2_delta_est = compracion_densidades_grafica(t = t2, df = base_filtrada_cauchy, var = delta_
  est); grafica2_delta_est
488
489 # Nota: Super interesante ver como cambia la distribuci n Cauchy
490 #   bajo los mismos para etros de localizaci n y escala
491 #   La distribuci n Cauchy es una distribuci n patol gica
492 #   en el sentido que su media y varianza es infinita
493
494 # plot(density(rcauchy(1000, location = 0, scale = 1)))
495
496 # Ejercicio e. ----
497
498 # I. ----
499
500 # Tama os de muestra
501 n20 = 20
502 n100 = 100
503 n200 = 200
504 n500 = 500
505 n1000 = 1000
506 n10000 = 10000
507
508 hetero20_no_robustos = delta(t = 1000, n = 20, distro = "normal", hetero = T, varianza_y0 = 2,
  varianza_y1 = 1, robustos = F)

```

```

509 hetero100_no_robustos = delta(t = 1000, n = 100, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = F)
510 hetero200_no_robustos = delta(t = 1000, n = 200, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = F)
511 hetero500_no_robustos = delta(t = 1000, n = 500, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = F)
512 hetero1000_no_robustos = delta(t = 1000, n = 1000, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = F)
513 hetero10000_no_robustos = delta(t = 1000, n = 10000, distro = "normal", hetero = T, varianza_y0 =
    2, varianza_y1 = 1, robustos = F)
514
515 n_vect_hetero = c(n20, n100, n200, n500, n1000, n10000)
516
517 # Dataframe con el tama o de muestra y la media y varianza del
518 # estimador de delta para los diferentes tama os de muestra
519 media_varianza_hetero = media_varianza_delta(t = 1000, n_vect_hetero, distro = "normal", hetero = T
    , varianza_y0 = 2, varianza_y1 = 1, robustos = F); glimpse(media_varianza_hetero)
520
521 # Nota: Se observa que a pesar de la heterocedasticidad en Y_{i} hay insesgidez
522 #     En la estimaci n de delta
523 #     De igual forma, la varianza muestral tambi n disminuye a medida que aumenta
524 #     la muestra (independiente de que haya heterocedasticidad en Y_{i})
525
526 # Construir una funci n que me compute los intervalos de confianza
527 # Intervalos de confianza: Ya sea para intervalos cl sicos o para intervalos robustos
528 #     a la heterocedasticidad como los calculados por la matriz de varianzas y
    covarianzas Huber-White
529
530 inter_confianza = function(df_estimaciones, delta_true, int_conf){
531   norm_inf = qnorm((1 - int_conf)/2)
532   norm_sup = qnorm((1 - int_conf)/2 + int_conf)
533   lim_inf = df_estimaciones$delta_est + norm_inf * sqrt(df_estimaciones$var_ols)
534   lim_sup = df_estimaciones$delta_est + norm_sup * sqrt(df_estimaciones$var_ols)
535   contiene_o_no = c()
536   for (i in 1:length(lim_inf)){
537     if ((lim_inf[i] < delta_true) && (delta_true < lim_sup[i])){
538       contiene_o_no = append(contiene_o_no, 1)
539     }else{
540       contiene_o_no = append(contiene_o_no, 0)
541     }
542   }
543   # Variable del dataframe
544   df = data.frame(lim_inf, lim_sup, contiene_o_no)
545   return(df)
546 }
547
548 # Intervalos de confianza del 90 %
549 int_conf1000_no_robustos_90 = inter_confianza(hetero1000_no_robustos, delta_true = 3, int_conf =
    0.9)
550
551 # Intervalos de confianza del 95 %
552 int_conf1000_no_robustos_95 = inter_confianza(hetero1000_no_robustos, delta_true = 3, int_conf =
    0.95)
553
554 # Intervalos de confianza del 99 %
555 int_conf1000_no_robustos_99 = inter_confianza(hetero1000_no_robustos, delta_true = 3, int_conf =
    0.99)
556
557 # II. ----
558
559 # Porcentaje de los intervalos de confianza del 90 % que contienen el par metro verdadero
560 porcentaje_no_robustos_90 = sum(int_conf1000_no_robustos_90$contiene_o_no)/nrow(int_conf1000_no_
    robustos_90) * 100; porcentaje_no_robustos_90
561
562 # Porcentaje de los intervalos de confianza del 95 % que contienen el par metro verdadero
563 porcentaje_no_robustos_95 = sum(int_conf1000_no_robustos_95$contiene_o_no)/nrow(int_conf1000_no_
    robustos_95) * 100; porcentaje_no_robustos_95
564
565 # Porcentaje de los intervalos de confianza del 99 % que contienen el par metro verdadero
566 porcentaje_no_robustos_99 = sum(int_conf1000_no_robustos_99$contiene_o_no)/nrow(int_conf1000_no_
    robustos_99) * 100; porcentaje_no_robustos_99
567
568 # III. ----
569
570 hetero20_robustos = delta(t = 1000, n = 20, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = T)
571 hetero100_robustos = delta(t = 1000, n = 100, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = T)
572 hetero200_robustos = delta(t = 1000, n = 200, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = T)
573 hetero500_robustos = delta(t = 1000, n = 500, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = T)
574 hetero1000_robustos = delta(t = 1000, n = 1000, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = T)
575 hetero10000_robustos = delta(t = 1000, n = 10000, distro = "normal", hetero = T, varianza_y0 = 2,
    varianza_y1 = 1, robustos = T)
576
577 n_vect_hetero = c(n20, n100, n200, n500, n1000, n10000)
578
579 # Dataframe con el tama o de muestra y la media y varianza del
580 # estimador de delta para los diferentes tama os de muestra

```

```

581 media_varianza_hetero = media_varianza_delta(t = 1000, n_vect_hetero, distro = "normal", hetero = T
    , varianza_y0 = 2, varianza_y1 = 1, robustos = T); glimpse(media_varianza_hetero)
582
583 # III. C lculo de intervalos de confianza cuando se tienen errores robustos ----
584
585 # Intervalos de confianza del 90 %
586 int_conf1000_robustos_90 = inter_confianza(hetero1000_robustos, delta_true = 3, int_conf = 0.9)
587
588 # Intervalos de confianza del 95 %
589 int_conf1000_robustos_95 = inter_confianza(hetero1000_robustos, delta_true = 3, int_conf = 0.95)
590
591 # Intervalos de confianza del 99 %
592 int_conf1000_robustos_99 = inter_confianza(hetero1000_robustos, delta_true = 3, int_conf = 0.99)
593
594 # III. Porcentaje de los intervalos de confianza cuando se tienen errores robustos ----
595
596 # Porcentaje de los intervalos de confianza del 90 % que contienen el par metro verdadero
597 porcentaje_robustos_90 = sum(int_conf1000_robustos_90$contiene_o_no)/nrow(int_conf1000_robustos_90)
    * 100; porcentaje_robustos_90
598
599 # Porcentaje de los intervalos de confianza del 95 % que contienen el par metro verdadero
600 porcentaje_robustos_95 = sum(int_conf1000_robustos_95$contiene_o_no)/nrow(int_conf1000_robustos_95)
    * 100; porcentaje_robustos_95
601
602 # Porcentaje de los intervalos de confianza del 99 % que contienen el par metro verdadero
603 porcentaje_robustos_99 = sum(int_conf1000_robustos_99$contiene_o_no)/nrow(int_conf1000_robustos_99)
    * 100; porcentaje_robustos_99
604
605 # Conclusi n: Los errores robustos c lculados con matriz de varianzas y covarianzas Huber-White,
606 # s sirven cu ndo hay heterocedasticidad en los errores de modelo
607 # Por el contrario, si se utilizan intervalos de confianza cl sicos
608 # sin corregir por heterocedasticidad, se observa que los intervalos
609 # de confianza cl sicos de manera sistem tica para los niveles de
610 # significancia sobreestiman, es decir calculan intervalos de
611 # confianza m s amplios de los que en realidad deberian computarse
612 # dada la presencia de la heterocedasticidad

```