

EXAMEN EXTRAORDINARIO DE ECONOMETRÍA

Conteste cada pregunta en un cuadernillo diferente en dos horas y media.

Todas las cuestiones (a), (b), etc., tienen la misma puntuación.

1. Un modelo para estimar los efectos del tabaquismo en el ingreso anual (quizás a través de ausencias laborales debido a enfermedades o efectos sobre la productividad) es

$$\log(\text{income}) = \beta_0 + \beta_1 \text{cigs} + \beta_2 \text{educ} + \beta_3 \text{age} + \beta_4 \text{age}^2 + u_1 \quad (1)$$

donde *cigs* es el número de cigarrillos que se consumen al día, en promedio, y *educ* (años de educación) y *age* (edad) se suponen exógenas.

Para reflejar el hecho de que el consumo de cigarrillos podría estar determinado conjuntamente con los ingresos, se especifica una ecuación de demanda de cigarrillos,

$$\text{cigs} = \gamma_0 + \gamma_1 \log(\text{income}) + \gamma_2 \text{educ} + \gamma_3 \text{age} + \gamma_4 \text{age}^2 + \gamma_5 \log(\text{cigpric}) + \gamma_6 \text{restaurn} + u_2, \quad (2)$$

donde *cigpric* es el precio de una cajetilla de cigarrillos (en centavos) y *restaurn* es una variable binaria igual a la unidad si la persona vive en un estado cuyos restaurantes restringen el consumo de tabaco, y se supone que éstas son variables exógenas al individuo.

- (a) ¿Cómo se interpretan β_1 y γ_1 ? ¿Qué signos se esperarían para γ_5 y γ_6 ? ¿Bajo qué supuestos estarían las ecuaciones (1) y (2) identificadas, si es posible que lo estén?
- (b) Realice los contrastes que sean factibles para comprobar las dos condiciones para la identificación de las ecuaciones (1) y (2) dada la información en la Tabla 1. Comente cualquier supuesto adicional que se utilice.
- (c) Compare la estimación de la ecuación de ingreso mediante MCO y mediante MC2E. ¿Qué conclusiones puede extraer? Proporcione una estimación de la variación esperada en los ingresos de un aumento del consumo de tabaco en un cigarrillo al día. ¿Podría realizar un contraste válido con la información proporcionada para determinar si la variación es negativa?

Table 1: Tabla de Regresiones

Var. dependiente:	(1)	(2)	(3)	(4)	(5)	(6)
cigs	0.00173 (0.00143)					
educ	0.0604 (0.00745)	-0.450 (0.156)	-0.472 (0.156)	0.0397 (0.0115)	-0.000301 (0.0103)	-1.15e-10 (0.0103)
age	0.0577 (0.00920)	0.823 (0.136)	0.824 (0.137)	0.0938 (0.0172)	-0.000269 (0.0115)	-3.60e-10 (0.0116)
age ²	-0.000631 (0.0000984)	-0.00959 (0.00144)	-0.00958 (0.00146)	-0.00105 (0.000197)	0.00000214 (0.000123)	3.17e-12 (0.000124)
log(<i>cigpric</i>)		-0.351 (6.027)			0.439 (0.392)	
restaurn		-2.736 (1.001)			-0.0145 (0.0675)	
hatcigs				-0.0421 (0.0188)		
Constante	7.795 (0.208)	1.580 (25.19)	-0.332 (3.226)	7.781 (0.208)	-1.784 (1.668)	1.05e-08 (0.259)
Observaciones	807	807	807	807	807	807
R ²	0.165	0.071	0.044	0.169	0.002	0.001

Errores estándar en paréntesis. Todas las regresiones están ajustadas por MCO.

hatcigs son las predicciones de cigs en la regresión de la columna (2).

hatu2 son los residuos del ajuste MC2E de la función de demanda, ecuación (2), usando *lcigpric* y *restaurn* como instrumentos para cigs

2. Se desea estimar un modelo de probabilidad lineal para estudiar los determinantes que llevan a la suscripción de un plan de pensiones,

$$pp = \beta_0 + \beta_1 \log(\text{income}) + \beta_2 \text{age} + \beta_3 \text{age}^2 + \beta_4 \text{male} + \beta_5 \text{married} + \beta_6 \text{male} * \text{married} + u, \quad (3)$$

donde $pp = 1$ si el individuo ha suscrito un plan y 0 en caso contrario, *income* es el ingreso anual, *age* la edad en años, *married* es una variable binaria indicando si el individuo está casado y *male* una variable binaria indicando si el individuo es un hombre. El modelo (3) se ha estimado con una muestra de 9275 individuos y los resultados están resumidos en la Tabla 2, junto con los de otros modelos relacionados.

- Interprete el coeficiente β_1 en (3) y proporcione un intervalo de confianza al 95% para el efecto sobre pp de un aumento de los ingresos del 10%.
- Determine si la probabilidad de suscribir un plan de pensiones de un hombre soltero es diferente a la de una mujer soltera usando el modelo (3). ¿Cuál es la diferencia estimada entre dicha probabilidad para un hombre y una mujer, ambos casados?
- Explique por qué con la información proporcionada no es posible llevar a cabo un contraste de hipótesis válido acerca de si la probabilidad de suscribir un plan de pensiones depende del género y estado civil del individuo.
- Finalmente, se decide estimar modelos diferentes para hombres y para mujeres incluyendo como variables explicativas $\log(\text{income})$, age , age^2 y *married*. Proporcione la ecuación estimada para los hombres usando los resultados de la columna (3).

Table 2: Tabla de Regresiones

Variable dependiente:	(1)	(2)	(3)
	pp	pp	pp
$\log(\text{income})$	0.229 (0.00838)	0.217 (0.00765)	0.231 (0.00948)
age	0.0137 (0.00347)	0.0142 (0.00346)	0.0158 (0.00394)
age ²	-0.000160 (0.0000398)	-0.000165 (0.0000397)	-0.000181 (0.0000448)
male	-0.0151 (0.0143)		0.197 (0.177)
married	-0.0360 (0.0116)		-0.0374 (0.0118)
male*married	-0.00875 (0.0244)		-0.00145 (0.0258)
$\log(\text{income}) * \text{male}$			-0.0112 (0.0203)
age*male			-0.00761 (0.00867)
age ² *male			0.0000747 (0.000102)
Constante	-0.774 (0.0723)	-0.770 (0.0718)	-0.830 (0.0824)
Observaciones	9275	9275	9275
R^2	0.084	0.083	0.084

Errores estándar robustos en paréntesis

Todas las regresiones están ajustadas por MCO

3. Suponga que el modelo verdadero que se desea estimar para explicar los resultados de un examen es

$$Testscore_i = \beta_0 + \beta_1 STR_i + \beta_2 Effort_i + u_i$$

donde STR es una variable continua que aumenta cuando el tamaño de la clase aumenta. $Effort$ es una variable continua que aumenta cuando el estudiante pone más esfuerzo. Finalmente, el error u incluye todos los otros factores que se supone están incorrelados con las otras variables explicativas del modelo. Suponga que $\beta_1 > 0$ y $\beta_2 > 0$, $Cov(STR, Effort) > 0$, $Cov(STR, u) = 0$, $Cov(Effort, u) = 0$.

Sin embargo, el econométra estima la siguiente especificación,

$$Testscore_i = \gamma_0 + \gamma_1 STR_i + e_i. \quad (4)$$

- (a) Discuta si la estimación MCO de γ_1 en (4) está sesgada o no para β_1 , y si es así, cuál es el signo del sesgo.
- (b) Imagine que en lugar de $Effort$ se observa un indicador de la asistencia a clase, CA_i , para el que se sabe que $Cov(CA, Effort) > 0$, por lo que se cumple que

$$Effort_i = \alpha_0 + \alpha_1 CA_i + v_i, \quad Cov(CA, v) = 0, \alpha_1 > 0.$$

También se sabe que

$$Cov(STR, v) = 0.$$

Interprete esta condición.

- (c) ¿Cuál sería el signo esperado de δ_3 en el siguiente modelo de regresión?

$$Testscore_i = \delta_0 + \delta_1 STR_i + \delta_3 CA_i + w_i$$

¿Cree que el estimador MCO de δ_1 sería sesgado o insesgado para β_1 ? ¿Por qué?

ALGUNOS VALORES CRÍTICOS: $Z_{0.90} = 1.282$, $Z_{0.95} = 1.645$, $Z_{0.975} = 1.96$, $\chi_{2,0.95}^2 = 5.99$, $\chi_{2,0.975}^2 = 7.378$, $\chi_{3,0.95}^2 = 7.815$, $\chi_{3,0.975}^2 = 9.348$, $\chi_{4,0.95}^2 = 9.488$, $\chi_{4,0.975}^2 = 11.143$, $\chi_{5,0.95}^2 = 11.071$, $\chi_{5,0.975}^2 = 12.833$, $\chi_{6,0.95}^2 = 12.592$, $\chi_{6,0.975}^2 = 14.449$, donde $\mathbb{P}(Z \leq Z_\alpha) = \alpha$ y $\mathbb{P}(\chi_m^2 \leq \chi_{m,\alpha}^2) = \alpha$, Z está distribuida como una normal estándar con media cero y varianza uno, y χ_m^2 como una chi-cuadrado con m grados de libertad.

EXAMEN EXTRAORDINARIO DE ECONOMETRÍA SOLUCIONES

1. Un modelo para estimar los efectos del tabaquismo en el ingreso anual (quizás a través de ausencias laborales debido a enfermedades o efectos sobre la productividad) es

$$\log(\text{income}) = \beta_0 + \beta_1 \text{cigs} + \beta_2 \text{educ} + \beta_3 \text{age} + \beta_4 \text{age}^2 + u_1 \quad (5)$$

donde *cigs* es el número de cigarrillos que se consumen al día, en promedio, y *educ* (años de educación) y *age* (edad) se suponen exógenas.

Para reflejar el hecho de que el consumo de cigarrillos podría estar determinado conjuntamente con los ingresos, se especifica una ecuación de demanda de cigarrillos,

$$\text{cigs} = \gamma_0 + \gamma_1 \log(\text{income}) + \gamma_2 \text{educ} + \gamma_3 \text{age} + \gamma_4 \text{age}^2 + \gamma_5 \log(\text{cigpric}) + \gamma_6 \text{restaurn} + u_2, \quad (6)$$

donde *cigpric* es el precio de una cajetilla de cigarrillos (en centavos) y *restaurn* es una variable binaria igual a la unidad si la persona vive en un estado cuyos restaurantes restringen el consumo de tabaco, y se supone que éstas son variables exógenas al individuo.

- (a) ¿Cómo se interpretan β_1 y γ_1 ?

[15%] β_1 : si se consume un cigarrillo más al día, el ingreso anual se incrementa (aproximadamente) en un $100\beta_1\%$ en promedio, todo lo demás igual.

[15%] γ_1 : si se incrementa el ingreso en un 1%, la demanda de cigarrillos aumenta (aproximadamente) en $\gamma_1/100$ unidades (cigarrillos) en promedio, todo lo demás igual.

[20%] ¿Qué signos se esperarían para γ_5 y γ_6 ?

γ_5 debería ser negativo (si aumenta el precio, la demanda de cigarrillos baja) y γ_6 debería ser negativo (si hay restricciones al consumo, la demanda de cigarrillos debería bajar, todo lo demás igual).

¿Bajo qué supuestos estarían las ecuaciones (5) y (6) identificadas, si es posible que lo estén?

[30%] La ecuación (5) está identificada si los dos posibles instrumentos, $\log(\text{cigpric})$ y *restaurn*, son exógenos en esa ecuación, es decir, se cumple que

$$H_0^{(exog)} : \text{Cov}(\log(\text{cigpric}), u_1) = 0 \text{ y } \text{Cov}(\text{restaurn}, u_1) = 0,$$

y si al menos uno de ellos es relevante, es decir, si en la forma reducida de *cigs*,

$$\text{cigs} = \pi_0 + \pi_1 \text{educ} + \pi_2 \text{age} + \pi_3 \text{age}^2 + \pi_4 \log(\text{cigpric}) + \pi_5 \text{restaurn} + v$$

la hipótesis nula

$$H_0^{(no\ rel)} : \pi_4 = \pi_5 = 0$$

es falsa, por lo que

$$H_1^{(no\ rel)} : \pi_4 \neq 0 \text{ y/o } \pi_5 \neq 0$$

es cierta.

[20%] La ecuación (6) no puede estar identificada porque no hay instrumentos exógenos que puedan usarse para el regresor endógeno $\log(\text{income})$.

- (b) Realice los contrastes que sean factibles para comprobar las dos condiciones para la identificación de las ecuaciones (5) y (6) dada la información en la Tabla 1. Comente cualquier supuesto adicional que se utilice.

[50%] Para contrastar $H_0^{(exog)}$ en contra de

$$H_1^{(exog)} : \text{Cov}(\log(\text{cigpric}), u_1) \neq 0 \text{ y/o } \text{Cov}(\text{restaurn}, u_1) \neq 0$$

podemos emplear el contraste de sobre-identificación de Hansen/Sargan, que se basa en el estadístico F para contrastar la hipótesis

$$H_0 : \delta_4 = \delta_5 = 0$$

en contra de

$$H_1 : \delta_4 \neq 0 \text{ y/o } \delta_5 \neq 0$$

en la regresión de los residuos MC2E sobre todas las variables exógenas,

$$\hat{u}_2 = \delta_0 + \delta_1 educ + \delta_2 age + \delta_3 age^2 + \delta_4 \log(cigpric) + \delta_5 restaurn + e.$$

Para ello podemos construir el estadístico F sólo válido para homocedasticidad (que es un supuesto adicional), comparando las columnas (5) = modelo no restringido y (6) = modelo restringido,

$$F = \frac{R_{no\ res}^2 - R_{res}^2}{1 - R_{no\ res}^2} \frac{n - k_{no\ res} - 1}{m} = \frac{0.002 - 0.001}{1 - 0.002} \frac{807 - 5 - 1}{2} = 0.4013$$

donde $k_{no\ res} = 5$ es el número de regresores en la ecuación no restringida y $m = 2$ es igual al número de restricciones (igual al número de instrumentos) y el estadístico J ,

$$J = mF = 2 \cdot 0.4013 = 0.8026.$$

Dado el número de regresores endógenos, $k = 1$, se compara J con el valor crítico de una $\chi_{m-k}^2 = \chi_1^2$, donde $m - k = 1$ es el grado de sobreidentificación, que al nivel de significación del 5% es igual a 3.86. Por lo tanto el estadístico J no es significativo al 5%, y no podemos rechazar la hipótesis nula de exogeneidad.

[50%] Para contrastar la hipótesis nula

$$H_0^{(no\ rel)} : \pi_4 = \pi_5 = 0$$

en contra de

$$H_1^{(no\ rel)} : \pi_4 \neq 0 \text{ y/o } \pi_5 \neq 0$$

usamos también un estadístico F sólo válido para homocedasticidad, comparando las columnas (2) = modelo no restringido y (3) = modelo restringido,

$$F = \frac{R_{no\ res}^2 - R_{res}^2}{1 - R_{no\ res}^2} \frac{n - k_{no\ res} - 1}{q} = \frac{0.071 - 0.044}{1 - 0.071} \frac{807 - 5 - 1}{2} = 11.640,$$

donde $k_{no\ res} = 5$ es el número de regresores en la ecuación no restringida y $q = 2$ es igual al número de restricciones (igual al número de instrumentos) y comparando el valor de F con el valor crítico de una $\chi_q^2/2 = \chi_2^2/2$, que al nivel de significación del 5% es igual a $5.99/2 \approx 3.00$, el estadístico F es significativo, rechazando $H_0^{(no\ rel)}$ y concluyendo que los instrumentos son relevantes, y además $F > 10$, por lo que podemos afirmar que los instrumentos no son débiles.

- (c) Compare la estimación de la ecuación de ingreso mediante MCO y mediante MC2E. ¿Qué conclusiones puede extraer?

[40%] La estimación MCO (columna 1) y la estimación MC2E (columna 4, segunda etapa de MC2E) de la ecuación de ingreso dan resultados muy diferentes: el coeficiente estimado para *cigs* cambia de signo (0.00173 y -0.0421 , respectivamente), por lo que podemos sospechar que la estimación MCO es inconsistente por un problema de endogeneidad de *cigs* en esa ecuación.

Proporcione una estimación de la variación esperada en los ingresos de un aumento del consumo de tabaco en un 1 cigarrillo al día.

[30%] Si $\Delta cigs = 1$, entonces los ingresos aumentarían en un $100 \cdot (-0.0421) \%$ en promedio, es decir bajarían en un 4,2%, usando la estimación MC2E.

¿Podría realizar un contraste válido con la información proporcionada para determinar si la variación es negativa?

[30%] No, porque los errores estándar de la estimación de la segunda etapa de MC2E obtenidos por MCO en la columna (4) reemplazando *cigs* por *haticgs* no son válidos al no tener en cuenta los efectos de la primera etapa (se deberían usar e.e. específicos).

2. Se desea estimar un modelo de probabilidad lineal para estudiar los determinantes que llevan a la subscripción de un plan de pensiones,

$$pp = \beta_0 + \beta_1 \log(\text{income}) + \beta_2 \text{age} + \beta_3 \text{age}^2 + \beta_4 \text{male} + \beta_5 \text{married} + \beta_6 \text{male} * \text{married} + u, \quad (7)$$

donde $pp = 1$ si el individuo ha suscrito un plan y 0 en caso contrario, income es el ingreso anual, age la edad en años, married es una variable binaria indicando si el individuo está casado y male una variable binaria indicando si el individuo es un hombre. El modelo (7) se ha estimado con una muestra de 9275 individuos y los resultados están resumidos en la Tabla 2, junto con los de otros modelos relacionados.

- (a) Interprete el coeficiente β_1 en (7) y proporcione un intervalo de confianza al 95% para el efecto sobre pp de un aumento de los ingresos del 10%.

[50%] Si el ingreso aumenta en un 1%, la probabilidad de que el individuo suscriba un plan de pensiones aumenta en aproximadamente en $\beta_1/100$ (o en β_1 puntos porcentuales si expresamos la probabilidad en tanto por cien) todo lo demás igual.

[50%] Intervalo de confianza al 95%, usando los resultados de la columna (1)

$$\begin{aligned} 10\hat{\beta}_1/100 \pm 1.96s.e. \left(10\hat{\beta}_1/100 \right) &= \hat{\beta}_1/10 \pm 1.96s.e. \left(\hat{\beta}_1 \right) /10 \\ 0.229/10 \pm 1.96 * 0.00838/10 &: [0.0213, 0.0245] \end{aligned}$$

- (b) Determine si la probabilidad de suscribir un plan de pensiones de un hombre soltero es diferente a la de una mujer soltera usando el modelo (7).

[50%] Comparando

$$\begin{aligned} \Pr(pp = 1 | \text{ hombre soltero}) &= \Pr(pp = 1 | \text{ male} = 1, \text{ married} = 0, \text{ age}, \text{ income}) \\ &= \beta_0 + \beta_1 \log(\text{income}) + \beta_2 \text{age} + \beta_3 \text{age}^2 + \beta_4 \end{aligned}$$

y

$$\begin{aligned} \Pr(pp = 1 | \text{ mujer soltera}) &= \Pr(pp = 1 | \text{ male} = 0, \text{ married} = 0, \text{ age}, \text{ income}) \\ &= \beta_0 + \beta_1 \log(\text{income}) + \beta_2 \text{age} + \beta_3 \text{age}^2 \end{aligned}$$

tomando los demás factores como fijos, el contraste buscado es

$$\begin{aligned} H_0 &: \beta_4 = 0 \\ H_1 &: \beta_4 \neq 0 \end{aligned}$$

que se realiza mediante un estadístico t usando los errores estándar robustos de los resultados de la columna (1)

$$t = \frac{\hat{\beta}_4}{e.e.(\hat{\beta}_4)} = \frac{-0.0151}{0.0143} = -1.0559$$

por lo que no podemos rechazar H_0 a los niveles habituales de significatividad: no se rechaza que las dos probabilidades son iguales.

¿Cuál es la diferencia estimada entre dicha probabilidad para un hombre y una mujer, ambos casados?

[50%] Las probabilidades tomando como fijos los demás factores son

$$\begin{aligned} \Pr(pp = 1 | \text{ hombre casado}) &= \Pr(pp = 1 | \text{ male} = 1, \text{ married} = 1, \text{ age}, \text{ income}) \\ &= \beta_0 + \beta_1 \log(\text{income}) + \beta_2 \text{age} + \beta_3 \text{age}^2 + \beta_4 + \beta_5 + \beta_6 \end{aligned}$$

y

$$\begin{aligned}\Pr(pp = 1 | \text{mujer casada}) &= \Pr(pp = 1 | \text{male} = 0, \text{married} = 1, \text{age}, \text{income}) \\ &= \beta_0 + \beta_1 \log(\text{income}) + \beta_2 \text{age} + \beta_3 \text{age}^2 + \beta_5,\end{aligned}$$

por lo que usando las estimaciones de la columna (1)

$$\begin{aligned}\widehat{\Pr}(pp = 1 | \text{hombre casado}) - \widehat{\Pr}(pp = 1 | \text{mujer casada}) &= \hat{\beta}_4 + \hat{\beta}_6 \\ &= -0.0151 - 0.00875 \\ &= -0.02385\end{aligned}$$

un hombre casado tiene una probabilidad 0.02385 menor de suscribir un plan que una mujer casada, todo lo demás igual.

- (c) *Explique por qué con la información proporcionada no es posible llevar a cabo un contraste de hipótesis válido acerca de si la probabilidad de suscribir un plan de pensiones depende del género y estado civil del individuo.*

Porque aunque se podría calcular el valor del estadístico F sólo válido para homocedasticidad para contrastar

$$H_0 : \beta_4 = \beta_5 = \beta_6 = 0$$

comparando los R^2 de las columnas (1) y (2), este contraste no es válido para el modelo lineal de probabilidad porque los errores de este modelo siempre son heterocedásticos y necesitaríamos en su lugar un estadístico F robusto para hacer un contraste válido.

- (d) *Finalmente, se decide estimar modelos diferentes para hombres y para mujeres incluyendo como variables explicativas $\log(\text{income})$, age , age^2 y married . Proporcione la ecuación estimada para los hombres usando los resultados de la columna (3).*

La probabilidad es

$$\begin{aligned}\widehat{\Pr}(pp = 1 | \text{hombre}) &= \widehat{\Pr}(pp = 1 | \text{male} = 1, \text{age}, \text{married}, \text{income}) \\ &= (-0.830 + 0.197) + (0.231 - 0.0112) \log(\text{income}) \\ &\quad + (0.0158 - 0.00761) \text{age} + (-0.000181 + 0.0000747) \text{age}^2 \\ &\quad + (-0.0374 - 0.00145) \text{married} \\ &= -0.633 + 0.2198 \log(\text{income}) \\ &\quad + 0.00819 \text{age} - 1.063 \times 10^{-4} \text{age}^2 - 0.03885 \text{married}\end{aligned}$$

3. Suponga que el modelo verdadero que se desea estimar para explicar los resultados de un examen es

$$\text{Testscore}_i = \beta_0 + \beta_1 \text{STR}_i + \beta_2 \text{Effort}_i + u_i$$

donde STR es una variable continua que aumenta cuando el tamaño de la clase aumenta. Effort es una variable continua que aumenta cuando el estudiante pone más esfuerzo. Finalmente, el error u incluye todos los otros factores que se supone están incorrelados con las otras variables explicativas del modelo. Suponga que $\beta_1 > 0$ y $\beta_2 > 0$, $\text{Cov}(\text{STR}, \text{Effort}) > 0$, $\text{Cov}(\text{STR}, u) = 0$, $\text{Cov}(\text{Effort}, u) = 0$.

Sin embargo, el econométra estima la siguiente especificación,

$$\text{Testscore}_i = \gamma_0 + \gamma_1 \text{STR}_i + e_i. \quad (8)$$

- (a) Discuta si la estimación MCO de γ_1 en (8) está sesgada o no para β_1 , y si es así, cuál es el signo del sesgo.

[50%] Se debe comprobar si se cumplen las dos condiciones para sesgo por Variables Omitidas (VO):

- $\text{Cov}(\text{STR}, \text{Effort}) \neq 0$ (en realidad > 0).
- $\text{Cov}(\text{Effort}, e) \neq 0$, que es cierta porque comparando la ecuación estimada (8) con la verdadera tenemos que

$$e_i = \beta_2 \text{Effort}_i + u_i$$

y podemos calcular $\text{Cov}(\text{Effort}, e) = \beta_2 \text{Var}(\text{Effort}) > 0$ porque $\beta_2 > 0$, y por tanto la estimación MCO de γ_1 estará sesgada para β_1 .

También podemos comprobar que

$$\begin{aligned} \text{Cov}(e_i, \text{STR}_i) &= \text{Cov}(\beta_2 \text{Effort}_i + u_i, \text{STR}_i) \\ &= \beta_2 \text{Cov}(\text{Effort}_i, \text{STR}_i) + \text{Cov}(u_i, \text{STR}_i) \\ &= \beta_2 \text{Cov}(\text{Effort}_i, \text{STR}_i) > 0, \end{aligned}$$

por lo que no se cumple el supuesto 1 de MCO.

[50%] El sesgo de estimación es

$$\text{Sesgo}(\hat{\gamma}_1) = E(\hat{\gamma}_1 | \mathbf{X}) - \beta_1 \approx \frac{\text{Cov}(e_i, \text{STR}_i)}{\text{Var}(\text{STR}_i)} = \beta_2 \frac{\text{Cov}(\text{Effort}_i, \text{STR}_i)}{\text{Var}(\text{STR}_i)} > 0,$$

es decir $\hat{\gamma}_1$ sobreestima sistemáticamente el valor verdadero de β_1 .

- (b) Imagine que en lugar de Effort se observa un indicador de la asistencia a clase, CA_i , para el que se sabe que $\text{Cov}(CA, \text{Effort}) > 0$, por lo que se cumple que

$$\text{Effort}_i = \alpha_0 + \alpha_1 CA_i + v_i, \quad \text{Cov}(CA, v) = 0, \alpha_1 > 0.$$

También se sabe que

$$\text{Cov}(\text{STR}, v) = 0.$$

Interprete esta condición.

La incorrelación entre STR y v indica que la parte de Effort que no puede explicar CA , tampoco está correlada con STR , es decir que STR no es una variable omitida en la regresión de Effort sobre CA , y el error de reemplazar Effort por CA no estará relacionado con STR (lo que convierte a CA en una variable de control válida para reemplazar la variable omitida Effort y poder estimar el valor verdadero del coeficiente de STR).

(c) ¿Cuál sería el signo esperado de δ_3 en el siguiente modelo de regresión?

$$Testscore_i = \delta_0 + \delta_1 STR_i + \delta_3 CA_i + w_i$$

[30%] En esta regresión el coeficiente δ_3 va a capturar la correlación del error $e_i = \beta_2 Effort_i + u_i$ con CA_i , por lo que esperaríamos que su signo sea el de $Cov(\beta_2 Effort_i, CA_i) = \beta_2 Cov(Effort_i, CA_i) > 0$, ya que en principio CA no está omitida el modelo verdadero y por tanto $Cov(CA, u) = 0$.

¿Cree que el estimador MCO de δ_1 sería sesgado o insesgado para β_1 ? ¿Por qué?

[30%] Usando el modelo verdadero y sustituyendo la ecuación de $Effort$ en términos de CA ,

$$\begin{aligned} Testscore_i &= \beta_0 + \beta_1 STR_i + \beta_2 Effort_i + u_i \\ &= \beta_0 + \beta_1 STR_i + \beta_2 (\alpha_0 + \alpha_1 CA_i + v_i) + u_i \\ &= \underbrace{(\beta_0 + \beta_2 \alpha_0)}_{=\delta_0} + \underbrace{\beta_1}_{=\delta_1} STR_i + \underbrace{\beta_2 \alpha_1}_{=\delta_3} CA_i + \underbrace{(\beta_2 v_i + u_i)}_{=w_i}, \end{aligned}$$

[40%] donde el error $w_i = \beta_2 v_i + u_i$ está incorrelado con STR_i (porque $Cov(STR, u) = 0$ y $Cov(STR, v) = 0$) y también esperamos que esté incorrelado con CA_i porque $Cov(CA, v) = 0$ y porque podemos suponer que $Cov(CA, u) = 0$ ya que CA no es una variable omitida en el modelo verdadero. En ese caso el modelo satisface los supuestos del modelo de regresión múltiple con la correspondiente definición de los coeficientes, y el estimador MCO de δ_1 estimaría consistentemente sin sesgo el coeficiente β_1 de STR en el modelo original, mientras que $\delta_3 = \beta_2 \alpha_1 = (+)(+) > 0$ como habíamos anticipado.