

Sesgo, Varianza del estimador y Multicolinealidad

José Eduardo Valencia Espinosa

29/9/2022

Entorno

```
library(pacman)
p.load(data.table, fixest, magrittr, sandwich, mvtnorm, MASS, sanwich, lme4, haven, woodridge)
set.seed(123)
```

1. Sesgo por variables omitidas

```
# Load data
library(wooldridge)
data("bweight")
bweight$cigs2 <- bweight$cigs + bweight$cigs^2 adding a column with squares
# Run regressions
est_0 <- feols(bweight ~ cigs, data = bweight)
est_1 <- feols(bweight ~ cigs + motheduc, data = bweight)
est_2 <- feols(bweight ~ cigs + motheduc + faminc, data = bweight)
est_3 <- feols(bweight ~ cigs + cigs2 + motheduc + faminc, data = bweight)
# nicer tables
# Standard errors (using fixest package)
etable(est_0, est_1, est_2, est_3, se = "standard") # Under MLRS
```

	est_0 <chr>	est_1 <chr>	est_2 <chr>
Dependent Var:	bweight	bweight	bweight
(Intercept)	119.8*** (0.5723)	115.4*** (3.107)	116.8*** (3.138)
cigs	-0.5138*** (0.0905)	-0.4862*** (0.0926)	-0.4633*** (0.0927)
motheduc		0.3308 (0.2328)	0.0143 (0.2580)
faminc			0.0915** (0.0325)
cigs2			
SE type	IID	IID	IID
Observations	1,388	1,387	1,387
1-10 of 12 rows 1-4 of 5 columns			Previous 2 Next

a. Usando est_2 prediga el valor en onzas del peso de un niño cuya madre fumó 2 cigarros por día durante el embarazo, con todo lo demás en su promedio.

```
values1 <- c(1, 2, mean(bweight$motheduc, na.rm=T), mean(bweight$faminc)) # omitiendo la fila con valores NA
coeff1 <- t(est_2$coefficients)
oz1 <- sum(coeff1*values1)
oz1
```

```
## [1] 118.7477
```

El modelo predice que el peso del bebé al nacer será de 118.7477 onzas (con 2 cigarros diarios durante el embarazo y todos los demás valores en su promedio).

b. ¿Cuál es el signo de la covarianza entre educación de la madre y número de cigarros fumados al día durante el embarazo? ¿Es lo que esperaba? Explique. (Vea Tabla 3.2. en Wooldridge ed.7)

```
cov(na.omit(bweight)$motheduc, na.omit(bweight)$cigs)
```

```
## [1] -2.799861
```

El signo de la covarianza es negativo, esto quiere decir que su relación es inversa. La relación anterior puede tener sentido debido a factores socioeconómicos, por ejemplo, que una mujer embarazada con mayor educación conozca los efectos negativos en la salud de su bebé si fuma durante el embarazo. Sin embargo, esta solo es una afirmación preliminar e intuitiva.

c. ¿Cuál es el efecto de fumar un cigarro extra en el embarazo en el peso en onzas considerando una relación no lineal? ¿Por qué utilizaría una relación no lineal? Explique.

Considerando est_3

```
est_3$coefficients[2:3]
```

```
##          cigs          cigs2
## -0.77874428  0.01229368
```

```
dw_bweight <- beta_0 + 2*beta_cigs(cigs) = -0.7787 + 0.0246 (cigs)
```

Utilizar una relación no lineal es útil para considerar que un cigarro más diario no afecta lo mismo que el anterior. Sin embargo, en este caso, esta relación no es estadísticamente significativa (al menos bajo el planteamiento de la regresión est.3).

d. Interprete R^2 y R^2 -adj en est_3. Intuitivamente ¿Por qué R^2 -adj es menor?

```
r1 <- data.frame(t(c(summary(lm(bweight ~ cigs + cigs2 + motheduc + faminc, data = bweight))$r.squared,
summary(lm(bweight ~ cigs + cigs2 + motheduc + faminc, data = bweight))$adj.r.squared)))
colnames(r1) <- c("r.squared", "adj.r.squared")
r1
```

	rsquared <dbl>	adj.rsquared <dbl>
	0.03173806	0.02893557
1 row		

Ambas R^2 s describen que tanto las variables del modelo explican las variaciones en el peso del bebé. Así, la regresión est.3 explica el 3.17% de las variaciones en el peso. La interpretación de la R ajustada es igual, sin embargo, está ajustada por el número de variables explicativas.

La R ajustada siempre será menor o igual a la R no ajustada pues la primera es una ponderación de la segunda en cuanto al número de regresores del modelo.

2. Varianza de los estimadores

```
# Standard errors (using fixest package)
# Table 1:
etable(est_0, est_1, est_2, est_3, se = "standard") # Under MLRS
# Table 2:
etable(est_0, est_1, est_2, est_3, se = "hetero") # Robust Standard Errors (HCl: Not White)
# Table 3:
etable(est_0, est_1, est_2, se = "cluster", cluster = c("white")) # One way Clustered
# Standard Errors using sandwich
# White's original SE (HCl) # This is with the classic White-Huber correction
etable(est_2, _vcov = sandwich::vcovHC, _vcov_args = list(type = "HC0"))
```

a. Demuestre matemáticamente por qué el SE en la segunda tabla es distinto al SE la Tabla 1, al menos para el estimador asociado a cigs.

La primera tabla asume homocedasticidad y la segunda ajusta la desviación estándar no constante por sus grados de libertad:

$$se(\hat{\beta}_j) = \frac{\sigma}{\sqrt{SST(1-R^2)}}, se(\hat{\beta}_j) = \frac{\hat{u}_j}{\sqrt{SST(1-R^2)} \cdot n-k-1}$$

b. Demuestre matemáticamente por qué SE en la tercera tabla es distinto de SE en la primera tabla. ¿Qué supuesto asumimos como vulnerado en la estimación de Tabla 3?

La tercera tabla hace el ajuste de White por clusters. Así asume que el supuesto de muestreo aleatorio fue vulnerado. Lo anterior porque asume que el muestreo aleatorio fue realizado de subgrupos y no de la población:

$$se(\hat{\beta}_j) = \frac{\hat{u}_j}{\sqrt{SST(1-R^2)}}$$

c. Dado: `cor(bweight$cigs, bweight$motheduc, use = "complete.obs")` y considerando est_0 (en ejercicio 1) ¿aumenta la varianza del estimador cigs al incluirse `motheduc` en est_1? Demuestre lo anterior utilizando las fórmulas de la varianza de beta para un modelo simple y para un múltiple.

```
cor(bweight$cigs, bweight$motheduc, use = "complete.obs")
```

```
## [1] -0.2138651
```

```
variance_comparison <- data.frame(t(c((est_0$se[2])^2, (est_1$se[2])^2)))
variance_comparison
```

	cigs <dbl>	cigs.1 <dbl>
	0.008188648	0.008579242
1 row		

Las varianzas son estadísticamente iguales.

d. ¿Cuál es el VIF Tras agregar `motheduc`? Discuta, considerando el trade-off sesgo varianza si debemos incluir `motheduc`.

```
VIF <- 1/(1-(cor(bweight$motheduc, bweight$cigs, use="complete"))^2)
VIF
```

```
## [1] 1.847931
```

Si debemos incluir `motheduc` en la regresión. En primer lugar, su VIF es pequeño (< 10). En segundo lugar, y en el mismo sentido, la varianza de `educ` no cambiaba al agregar `motheduc`, por lo que no añade "ruido" a nuestra regresión (las variables no son multicolineales). En tercer lugar, agregar esta variable de control ayuda a disminuir el sesgo de los β 's.

e. ¿A qué concepto nos referimos en el inciso c?

Multicolinealidad.

3. BEAUTY

```
beauty_age <- beauty$educ+beauty$exper+6
coefest(lm(beauty$lwage ~ beauty$looks + beauty$exper + beauty$exper*sq + beauty$educ + beauty_age),
vcov = vcovHC(lm(beauty$lwage ~ beauty$looks + beauty$exper + beauty$exper*sq + beauty$educ + beauty_age),
type = "HC1"))
```

```
##
## t test of coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.04627736  0.10839716  0.4269  0.669588
## beauty$looks  0.06286811  0.02218328  2.7980  0.005221 **
## beauty$exper  0.04873264  0.00453849 10.5489 < 2.2e-16 ***
## beauty$exper*sq -0.00098846  0.00030411  -6.7088  2.96e-11 ***
## beauty$educ   0.06836892  0.00593334 11.5228 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

a. Explore las variables con `beauty` y agregue una columna a los datos que incluya una aproximación de edad para cada `i`: `beauty$age <- beauty$educ+beauty$exper+6`

```
beauty$age <- beauty$educ+beauty$exper+6
```

b. Estime una regresión simple de `log(wage)` con `looks` como explicativa con errores estándar robustos (HC1): pegue el resultado abajo e interprete β_1

```
logwage1 <- lm(log(wage) ~ looks, data = beauty)
coefest(logwage1, type = "HC1")
```

```
##
## t test of coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.496483  0.079518 18.7978 < 2e-16 ***
## looks        0.058951  0.024432  2.8555  0.03723 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Por cada `rating` adicional en `apariencia` aumenta en 5.09% el salario.

c. Agregue como controles a la regresión en b. la experiencia, la experiencia al cuadrado, la educación en años y la edad ¿Qué pasa con su regresión y por qué? ¿Qué supuesto se puede estar violando?

```
exper2 <- beauty$exper^2
logwage2 <- lm(lm(wage) ~ looks + exper + (exper2) + educ + age, data = beauty)
coefest(logwage2, type = "HC1")
```

```
##
## t test of coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.04627735  0.10533299  0.4393  0.660488
## looks        0.06286811  0.02192319  2.8312  0.004712 **
## exper        0.04873264  0.00478836 10.1773 < 2.2e-16 ***
## exper2       -0.00098846  0.00010767  -9.4870 1.257e-10 ***
## educ         0.06836892  0.00581266 11.7821 < 2.2e-16 ***
## age          NA          NA          NA          NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

No calcula la β para la variable `age`. El problema es que esta variable es una combinación lineal de `educ` y de `exper`, es decir, existe colinealidad perfecta.

d. Solucione el problema en c. y vuelva a estimar su regresión. Pegue el resultado abajo e interprete β_1 .

```
logwage3 <- lm(log(wage) ~ looks + exper + (exper2) + age, data = beauty)
coefest(logwage3, type = "HC1")
```

```
##
## t test of coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.36393619  0.12920133 -2.8168  0.004926 **
## looks        0.06286811  0.02192318  2.8312  0.004712 **
## exper        -0.01963628  0.00788322  -2.4989  0.012871 *
## exper2       -0.00098846  0.00010767  -9.4870 1.257e-10 ***
## age         0.06836892  0.00581266 11.7821 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Eliminamos una de las variables colineales perfectamente. Controlando con las demás variables, ahora la β_1 cambia. Por cada `rating` adicional en `apariencia` aumenta en 6.2% el salario.

4. SESGO EN β

b. Omitting an important variable: sesga a la β ya que su covarianza con el error no es igual a 0.

5. SESGO EN VARIANZA

a. Heteroskedasticity: que la varianza no sea constante en toda la regresión la sesga.

b. A sample correlation coefficient of 0.95 between two independent variables both included in the model: por multicolinealidad (ya fue explicado).

6. PARTIALLING OUT

```
educ <- lm(log(wage) ~ exper + tenure, data=wage1)
r1 <- rsid(educ)
logwage4r1 <- lm(lwage ~ r1, data=wage1)
logwage5 <- lm(lwage ~ educ + exper + tenure, data=wage1)
```

a. Pegue abajo el resultado de las dos regresiones de interés y explique intuitivamente por qué se observa que el coeficiente de `r1` es igual al de `educ` cuando controla por `exper` y `tenure`.

```
beta_comparison <- c(summary(logwage4r1)$coefficients[2,1], summary(logwage5)$coefficients[2,1])
beta_comparison
```

```
## [1] 0.89202899 0.89202899
```

Lo anterior lo explica el teorema de Frisch-Waugh-Lovell. Intuitivamente, los residuos de la primera regresión lineal terminan de explicar la variación en educación. Así, al proyectar `educ` en la tercera regresión (controlando por las mismas variables `exper` y `tenure`), obtendrás el efecto de la educación en el salario.

7. HOW MANY STUDENTS

a. Apartir de la marca "From here". Explique cada línea pertinente del código.

El código establece un loop que realiza una regresión lineal 10,000 veces para dos distribuciones normales multivariadas. Lo que diferencia a estas dos es su covarianza (en `sigma`). Los coeficientes de cada regresión lineal los guarda en un vector vacío definido anteriormente (mediante `coordenados`).

b. Escriba la matriz de varianza – covarianza de los estimadores 1 y 2.

Para la estimación 1:

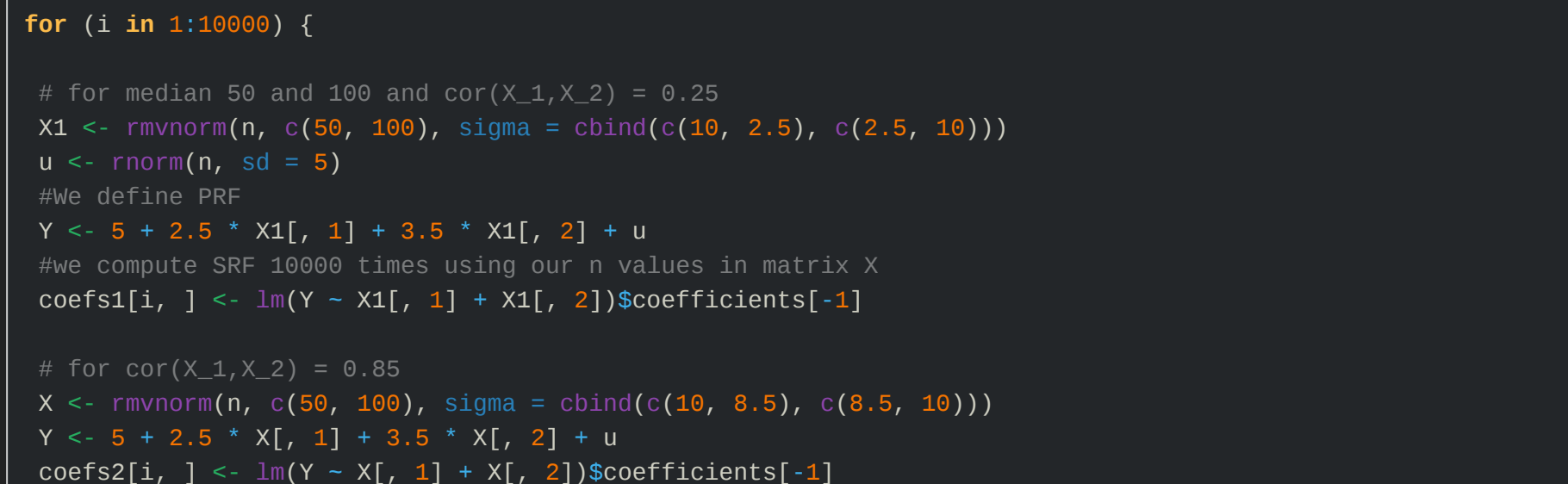
```
[10.3966129 3.2817823]
[3.2817823 9.9553062]
```

Para la estimación 2:

```
[10.4933833 9.2995875]
[9.2995875 11.0155659]
```

c. ¿Por qué los histogramas de `coef1` y `coef2` son distintos?

```
#Histogram of coefficients x1 and x2 for the two 10000s loops of regressions.
par(mfrow=c(1,2))
hist(coef1, ylim=c(0,5000))
hist(coef2, ylim=c(0,5000))
```



Porque las covarianzas son distintas y esto afecta a la varianza de los β 's

d. ¿Por qué las varianzas de `coef1` y `coef2` son distintas?

Porque entre más correlacionados (coeficiente de Pearson) estén dos regresores (cercano a `cor()` = ± 1), aumenta la varianza de los β 's (trade-off varianza y sesgo). "Aumenta el ruido"

e. ¿Qué parte del código genera estas diferencias?

En `sigma`, defina las covarianzas de manera distinta.

f. ¿Qué concepto/supuesto central estamos simulando en este ejercicio?

Multicolinealidad.

g. Ahora cambie `n` a 1000, ¿qué sucede con los histogramas y las varianzas? Explique intuitivamente, pero con el uso de alguna fórmula ¿por qué se da este cambio?

```
n <- 1000
coef1 <- cbind("beta_hat_1" = numeric(10000), "beta_hat_2" = numeric(10000))
coef2 <- coef1
##From here##
# Loop sampling and estimation (be patient it takes a while)
for (i in 1:10000) {
  # for median 50 and 100 and cov(X.1,X.2) = 0.25
  X1 <- rmvnorm(n, c(50, 100), sigma = cbind(c(10, 2.5), c(2.5, 10)))
  u <- rnorm(n, sd = 5)
  #we define PRF
  Y <- 5 + 2.5 * X1[, 1] + 3.5 * X1[, 2] + u
  #we compute SRP 10000 times using our n values in matrix X
  coef1[i, ] <- lm(Y ~ X1[, 1] + X1[, 2])$coefficients[-1]
  # for cov(X.1,X.2) = 0.85
  X <- rmvnorm(n, c(50, 100), sigma = cbind(c(10, 8.5), c(8.5, 10)))
  Y <- 5 + 2.5 * X[, 1] + 3.5 * X[, 2] + u
  coef2[i, ] <- lm(Y ~ X[, 1] + X[, 2])$coefficients[-1]
}
```

h. Aun bajo el supuesto MRL4 ¿Se puede obtener un estimador inconsistente en una muestra aleatoria de individuos? De un ejemplo con el uso de los histogramas.

Entre mayor sea la varianza de β , es más probable que los muestreos puedan ocurrir en los bordes de la distribución de β . De esta forma, aunque no esté sesgado, puede ser inconsistente.

Sin embargo, mientras mayor sea el número de observaciones, la $\hat{\beta}$ tenderá a β y este problema ya no será relevante:

$$\hat{\beta}_j = \beta_j + \frac{\text{cov}(u_j, x_j)}{\text{var}(x_j)} = \beta_j$$

Por ejemplo, dos histogramas de la distribución de β_1 y β_2 con correlación de 0.85 pero la primera con 100 observaciones y la segunda con 100,000 observaciones.

