

Análisis de varianza: heterocasticidad

José Eduardo València Espinosa

3/11/2022

1. Considera un modelo lineal que explica el consumo de cerveza mensual:

$$\begin{aligned} \text{beer} &= \beta_0 + \beta_1 \text{inc} + \beta_2 \text{price} + \beta_3 \text{educ} + \beta_4 \text{female} + u \\ E(u|\text{inc}, \text{price}, \text{educ}, \text{female}) &= 0 \\ \text{Var}(u|\text{inc}, \text{price}, \text{educ}, \text{female}) &= \sigma^2 \text{inc}^2 \end{aligned}$$

a. Escribe una ecuación transformada que tenga error homocástico. Muestra que es homocedástico.

$$\begin{aligned} \frac{\text{beer}}{\text{inc}} &= \frac{\beta_0}{\text{inc}} + \beta_1 + \beta_2 \frac{\text{price}}{\text{inc}} + \beta_3 \frac{\text{educ}}{\text{inc}} + \beta_4 \frac{\text{female}}{\text{inc}} + \frac{u}{\text{inc}} \\ \text{Var}\left(\frac{u}{\text{inc}}|\text{inc}, \text{price}, \text{educ}, \text{female}\right) &= \text{Var}(u|\text{inc}, \text{price}, \text{educ}, \text{female}) \cdot \frac{1}{\text{inc}^2} = \sigma^2 \end{aligned}$$

2. Usando los datos de GPA3, la siguiente ecuación fue estimada para los estudiantes de segundo semestre de Otoño:

$$\begin{aligned} \widehat{\text{trmgpa}} &= -2.12 & +0.900 \cdot \text{crsgpa} & -0.193 \cdot \text{cumgpa} & +0.0014 \cdot \text{tothrs} \\ & (0.55) & (0.165) & (0.064) & (0.0012) \\ & (0.55) & (0.166) & (0.074) & (0.0012) \\ & +0.0018 \cdot \text{sat} & -0.0039 \cdot \text{hsperc} & +0.351 \cdot \text{female} & -0.157 \cdot \text{season} \\ & (0.0002) & (0.0018) & (0.085) & (0.098) \\ & (0.0002) & (0.0019) & (0.079) & (0.080) \\ n &= 269, R^2 = 0.465 \end{aligned}$$

a. ¿Las variables crsgpa , cumgpa y tothrs tienen el efecto estimado esperado? ¿Cuál de estas variables es significativa al 5%? ¿Importa que errores estándar usas?

Para crsgpa , si esperas que mientras mayor haya sido el promedio de GPA de todos sus cursos, mayor será el GPA esperado de ese semestre (*ceteris paribus*). Lo mismo para cumgpa . En cuanto a tothrs , bien podría ser una variable irrelevante, pues cuántas horas acreditadas de semestres pasados. Lo mismo para un alumno no explica la capacidad que tiene de obtener un puntaje más alto o bajo. Por ejemplo, si ambos tienen el mismo crsgpa y cumgpa pero que un alumno haya tomado más cursos (probablemente por que sea de mayor semestre) no evidencia alguna capacidad adicional para un GPA más alto.

	E. No Rob.	E. Rob.	V. Crit.
t_{crsgpa}	5.1428	5.4216	±1.9688
t_{cumgpa}	3.0156	2.6081	±1.9688
t_{tothrs}	1.1666	1.1666	±1.9688

Por lo tanto, crsgpa y cumgpa son significativas al 5%. En este caso, la conclusión es la misma utilizando errores no robustos o robustos, por lo que no es importante. Sin embargo, puede que la n no sea lo suficientemente grande para utilizar errores robustos (por sus propiedades asintóticas). Si importara que errores estándar usas para la conclusión, entonces lo más apropiado sería obtener Mínimos Cuadrados Ponderados (si es que hay heterocasticidad).

b. ¿Por qué la hipótesis $H_0: \beta_{\text{crsgpa}} = 1$ hace sentido? Prueba esta hipótesis contra la alternativa de dos cosas al 5% de significancia (usando ambos errores estándar). Describe tus conclusiones.

Por que el promedio del GPA es el equivalente muestral al valor esperado. Así, manteniendo todo lo demás constante, el GPA del semestre se espera que sea el promedio de los cursos anteriores.

	E. No Rob.	E. Rob.	V. Crit.
$H_0: \beta_{\text{crsgpa}} = 1$			
$H_a: \beta_{\text{crsgpa}} \neq 1$			
t_{crsgpa}	-0.5714	-0.6024	±1.9688

No podemos rechazar la hipótesis nula al 5% de significancia. Este resultado va en la misma dirección que la justificación anterior (del promedio como estimador *ceteris paribus* del GPA).

c. Prueba si hay un efecto de temporada en el GPA del semestre (usando ambos errores estándar). ¿El nivel de significancia al que se puede rechazar la hipótesis nula depende de los errores estándar utilizados?

	E. No Rob.	E. Rob.	V.C. 1%	V.C. 5%	V.C. 10%
t_{season}	-1.6020	-1.9625	±2.5942	±1.9688	±1.6505

Si, si utiliza errores estándar no robustos, rechazamos la hipótesis nula para los 3 niveles de significancia. Si utiliza errores estándar robustos, rechazamos la hipótesis nula al 1 y 5% de significancia pero no la rechazamos al 10% de significancia.

3. Considera un modelo a nivel de empleado, donde la variable no observada f_i es un "efecto de la firma" de cada empleado de una firma dada. El término de error $v_{i,e}$ es específico al empleado e en la firma i . El error compuesto es $u_{i,e} = f_i + v_{i,e}$, como en la ecuación (8.28) Wooldridge 7ma ed.

i. Asume que $\text{Var}(f_i) = \sigma_f^2$, $\text{Var}(v_{i,e}) = \sigma_v^2$, y f_i y $v_{i,e}$ no correlacionados. Muestra que $\text{Var}(u_{i,e}) = \sigma_f^2 + \sigma_v^2$; llama a esto σ^2

$$\text{Var}(u_{i,e}) = \text{Var}(f_i + v_{i,e}) = \text{Var}(f_i) + \text{Var}(v_{i,e}) + 2\text{Cov}(f_i, v_{i,e})$$

Si no hay correlación entre f_i y $v_{i,e}$, entonces $\text{Cov}(f_i, v_{i,e}) = 0$ (f_i y $v_{i,e}$ son ortogonales). Esto es lo mismo que decir $\text{Cov}(f_i, v_{i,e}) = 0$. Entonces,

$$\text{Var}(u_{i,e}) = \text{Var}(f_i) + \text{Var}(v_{i,e}) + 2 \cdot 0 = \sigma_f^2 + \sigma_v^2$$

ii. Ahora, supón que para $e \neq g$, $v_{i,e}$ y $v_{i,g}$ no correlacionados. Muestra que $\text{Cov}(u_{i,e}, u_{i,g}) = \sigma_f^2$

$$u_{i,e} = f_i + v_{i,e}$$

$$u_{i,g} = f_i + v_{i,g}$$

$$\text{Cov}(u_{i,e}, u_{i,g}) = \text{Cov}(f_i + v_{i,e}, f_i + v_{i,g})$$

Por propiedades del producto punto,

$$\text{Cov}(f_i + v_{i,e}, f_i + v_{i,g}) = \text{Cov}(f_i, f_i) + \text{Cov}(v_{i,e}, f_i) + \text{Cov}(f_i, v_{i,g}) + \text{Cov}(v_{i,e}, v_{i,g})$$

Ya que $v_{i,e}$ y $v_{i,g}$ no correlacionados, $\text{Cov}(v_{i,e}, v_{i,g}) = 0$.

$$\text{Cov}(u_{i,e}, u_{i,g}) = \text{Var}(f_i) + \text{Cov}(v_{i,e}, f_i) + \text{Cov}(f_i, v_{i,g})$$

Con e y g no correlacionados, $\text{Cov}(v_{i,e}, f_i)$, $\text{Cov}(f_i, v_{i,g}) = 0$.

$$\text{Cov}(u_{i,e}, u_{i,g}) = \text{Var}(f_i) = \sigma_f^2$$

iii. Define $\bar{u}_i = m_i^{-1} \sum_{e=1}^{m_i} u_{i,e}$ como el promedio del error compuesto de la firma i . Muestra que $\text{Var}(\bar{u}_i) = \sigma_f^2 + \frac{\sigma_v^2}{m_i}$

$$\text{Var}(\bar{u}_i) = \text{Var}(m_i^{-1} \sum_{e=1}^{m_i} u_{i,e}) = \frac{1}{m_i^2} \text{Var}(\sum_{e=1}^{m_i} u_{i,e})$$

Sin correlación serial,

$$\text{Var}(\bar{u}_i) = \frac{1}{m_i^2} \sum_{e=1}^{m_i} \text{Var}(u_{i,e}) = \frac{1}{m_i^2} \sum_{e=1}^{m_i} \text{Var}(m_i f_i + v_{i,e})$$

Sin correlación entre f_i y $v_{i,e}$,

$$\text{Var}(\bar{u}_i) = \frac{1}{m_i^2} \sum_{e=1}^{m_i} (m_i \sigma_f^2 + \sigma_v^2) = \frac{1}{m_i^2} \cdot m_i \cdot (m_i \sigma_f^2 + \sigma_v^2) = \sigma_f^2 + \frac{\sigma_v^2}{m_i}$$

4. Use los datos de wage1

```
library(pacman)
p_load(wooldridge, ggplot2, ggthemes, lme4, fixest)
```

a. Use MCO para estimar un modelo que relacione el salario con la educación, experiencia y permanencia.

```
lmwage <- lm(wage ~ educ + exper + tenure, data=wage1)
```

b. Computa manualmente el test Breusch-Pagan, incluyendo al estadístico F y LM.

$$u^2 = \delta_0 + \delta_1 \cdot \text{educ} + \delta_2 \cdot \text{exper} + \delta_3 \cdot \text{tenure} + \text{error}$$

$$H_0: \delta_1 = \delta_2 = \delta_3 = 0$$

Así, el modelo restringido es $H_a: \hat{u}^2 = \delta_1$.

Con estadístico F:

$$F = \frac{R_{\text{edu}}^2/3}{(1 - R_{\text{edu}}^2)/526 - 3 - 1}$$

```
wage1$u_sqr <- resid(lmwage)^2
lmBP <- lm(u_sqr ~ educ + exper + tenure, data=wage1)
F_BP <- (sum(resid(lmwage)^2)/3) / ((1 - sum(resid(lmwage)^2)/526) / 3)
VC.F <- qf(0.05, 3, 526 - 4, lower.tail = F)
```

F_{BP} VC : 5%

15.5282 > 1.1547

∴ Se rechaza hipótesis nula. El test Breusch-Pagan con el estadístico F concluye heterocasticidad.

Con estadístico LM:

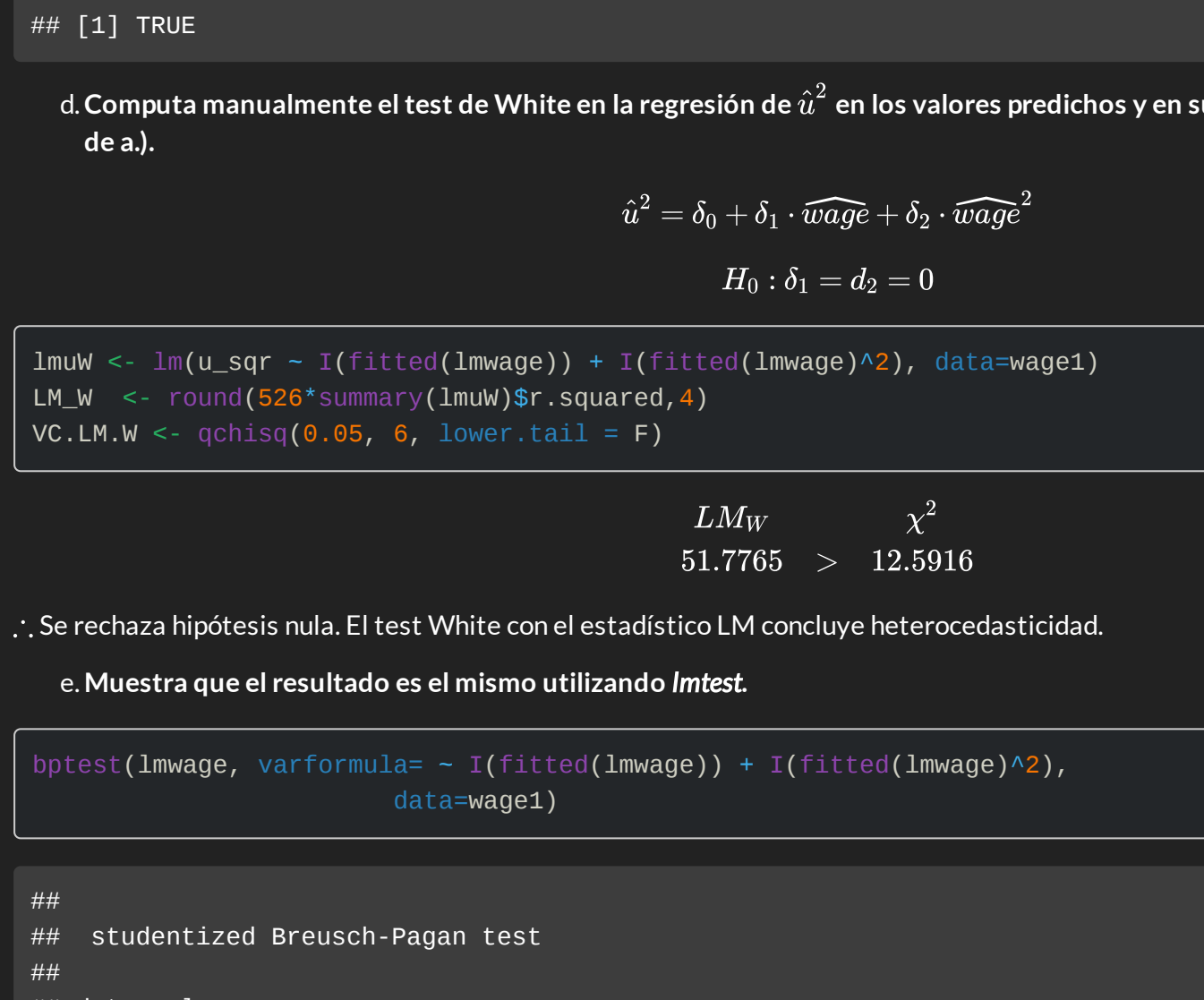
$$LM = 826 \cdot R_{\text{edu}}^2$$

```
LM_BP <- round(sum(resid(lmwage)^2)/3, 4)
VC.LM_BP <- qchisq(0.05, 3, lower.tail = F)
```

LM_{BP} χ^2

43.0956 > 7.8147

∴ Se rechaza hipótesis nula. El test Breusch-Pagan con el estadístico LM concluye heterocasticidad.



c. Muestra que lo que obtuviste en b. lo puedes obtener con la siguiente función:

```
bptest(lmwage)

##
## Studentized Breusch-Pagan test
##
## data: lmwage
## BP = 43.096, df = 3, p-value = 2.349e-09

BP.lmwage <- round(bptest(lmwage)$statistic["BP"], 4)
BP.lmwage == LM_BP

## [1] TRUE
```

d. Computa manualmente el test de White en la regresión de \hat{u}^2 en los valores predichos y en sus cuadrados (utilizando regresión de a.).

$$\hat{u}^2 = \delta_0 + \delta_1 \cdot \widehat{\text{wage}} + \delta_2 \cdot \widehat{\text{wage}}^2$$

$$H_0: \delta_1 = \delta_2 = 0$$

```
lmw <- lm(u_sqr ~ I(fitted(lmwage)) + I(fitted(lmwage)^2), data=wage1)
LM.W <- round(sum(resid(lmw)^2)/3, 4)
VC.LM.W <- qchisq(0.05, 2, lower.tail = F)
```

LM_W χ^2

51.7765 > 12.5916

∴ Se rechaza hipótesis nula. El test White con el estadístico LM concluye heterocasticidad.

e. Muestra a qué resultado es el mismo utilizando lme4.

```
bptest(lmwage, varformulae ~ I(fitted(lmwage)) + I(fitted(lmwage)^2),
      data=wage1)

##
## Studentized Breusch-Pagan test
##
## data: lmwage
## BP = 51.777, df = 2, p-value = 5.713e-12

lmwage <- round(bptest(lmwage, varformulae ~ I(fitted(lmwage)) + I(fitted(lmwage)^2),
      data=wage1)$statistic["BP"], 4)
lmwage == LM.W

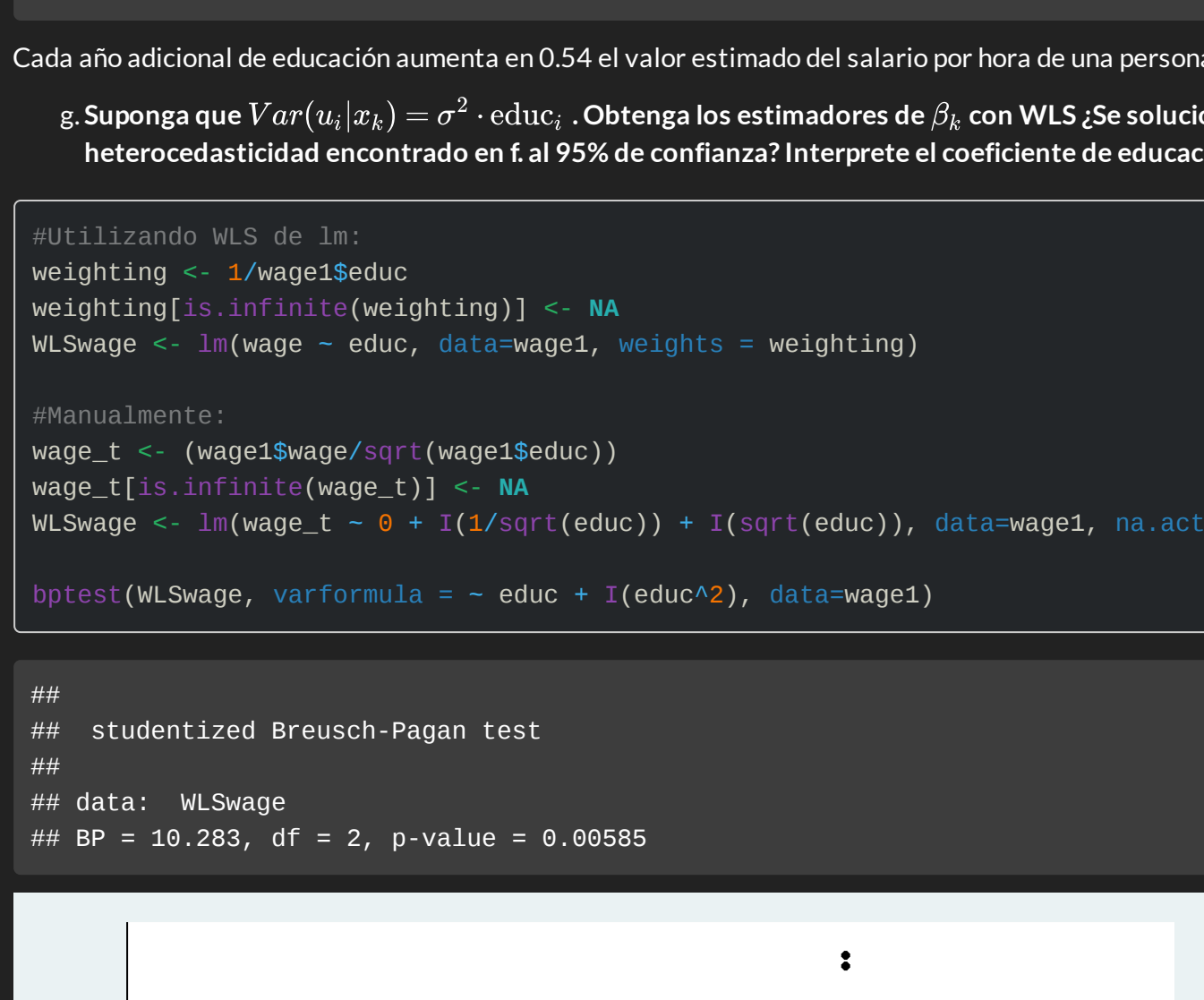
## [1] TRUE
```

f. Usa solo educación como predictor e interprete $\hat{\beta}$. ¿Es heterocedística la relación? Interprete el coeficiente de educ.

```
lmwage2 <- lm(wage ~ educ, data=wage1)
bptest(lmwage2, varformulae ~ educ + I(educ^2), data=wage1)

##
## Studentized Breusch-Pagan test
##
## data: lmwage2
## BP = 23.241, df = 2, p-value = 8.982e-06
```

Según el test de White, lmwage2 es heterocedástico.



```
lmwage2$coefficients

## (Intercept)      educ
## -0.8048516    0.5433593
```

Cada año adicional de educación aumenta en 0.54 el valor estimado del salario por hora de una persona.

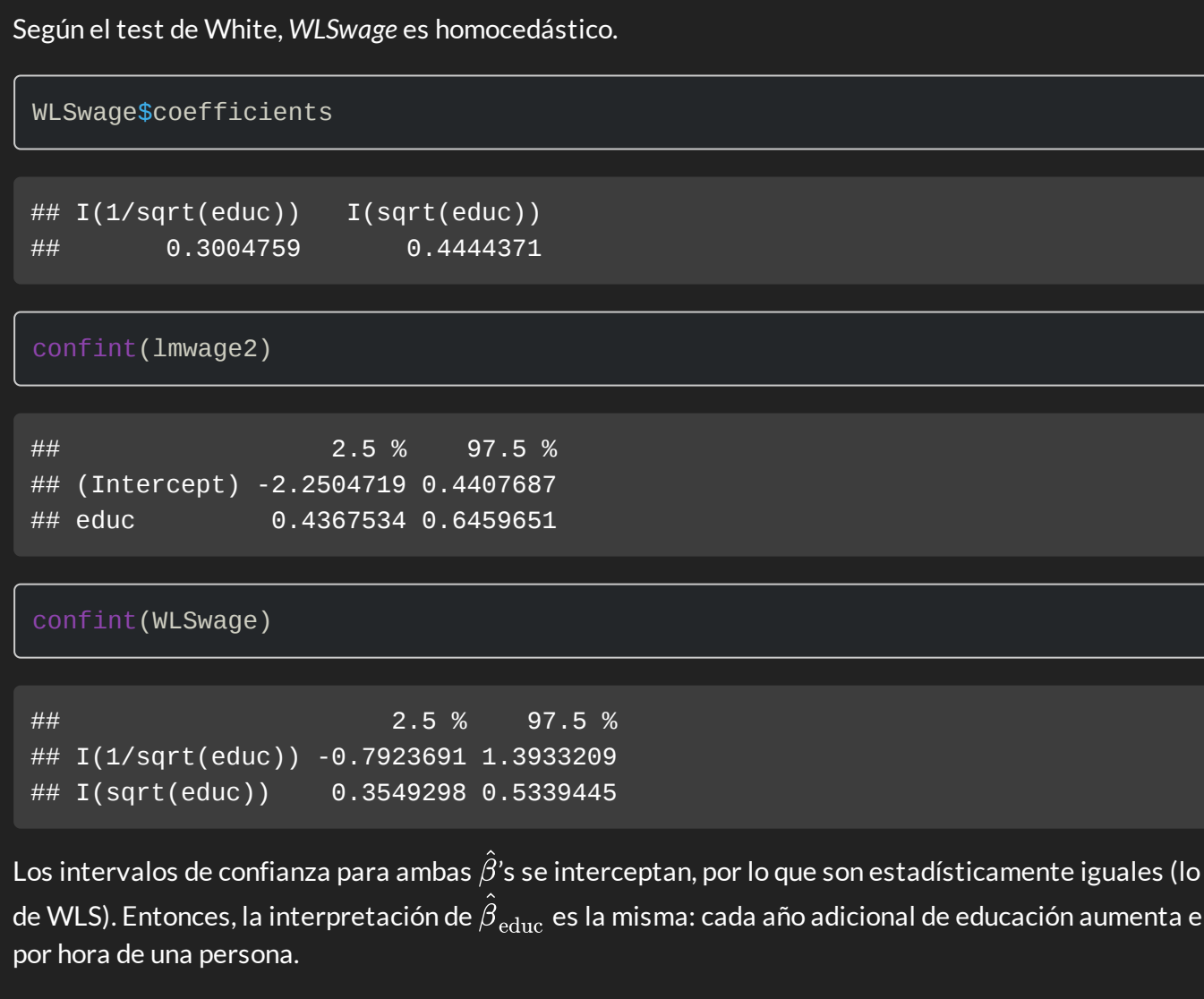
g. Suponga que $\text{Var}(u_i|x_i) = \sigma^2 \cdot \text{educ}_i$. Obtenga los estimadores de β_0 con WLS. ¿Se soluciona el problema de heterocedasticidad encontrado en f. al 95% de confianza? Interprete el coeficiente de educación.

```
#utilizando WLS de lm
weighting <- 1/wage1$educ
weighting[is.infinite(weighting)] <- NA
WLSwage <- lm(wage ~ educ, data=wage1, weights = weighting)

#Manualmente:
wage_t <- (wage1$wage/sqrt(wage1$educ))
wage_t[is.infinite(wage_t)] <- NA
lmw_t <- lm(wage_t ~ 0 + I(1/sqrt(educ)) + I(sqrt(educ)), data=wage1, na.action = na.exclude)
WLSwage <- lm(wage_t ~ 0 + I(1/sqrt(educ)) + I(sqrt(educ)), data=wage1)

bptest(WLSwage, varformulae ~ educ + I(educ^2), data=wage1)

##
## Studentized Breusch-Pagan test
##
## data: WLSwage
## BP = 10.283, df = 2, p-value = 0.00585
```



Según el test de White, WLSwage es homocedástico.

```
WLSwage$coefficients

## I(1/sqrt(educ)) I(sqrt(educ))
## 0.3904759    0.4444371

confint(lmwage2)

##               2.5 %      97.5 %
## (Intercept) -2.2504719 0.4487687
## educ         0.4367534 0.6459551

confint(WLSwage)

##               2.5 %      97.5 %
## I(1/sqrt(educ)) -0.7923591 1.3933289
## I(sqrt(educ))   0.3549298 0.5339445
```

Los intervalos de confianza para ambos $\hat{\beta}$'s se interceptan, por lo que son estadísticamente iguales (lo esperado según la transformación de WLS). Entonces, la interpretación de $\hat{\beta}_{\text{educ}}$ es la misma: cada año adicional de educación aumenta en 0.44 el valor estimado del salario por hora de una persona.

h. Verifica que los valores ajustados (fitted) en d. son todos positivos y obtén el WLS estimando h_i (FGLS).

```
#son positivos
sum(fitted(lmw)) < 0

## [1] 0

lmuh <- lm(log(u_sqr) ~ educ + exper + tenure, data= wage1)
h_i <- exp(fitted(lmuh))
#son positivos
sum(h_i < 0)

## [1] 0

#es con h_i estimado
FGLSwage <- lm(wage ~ educ + exper + tenure, data= wage1, weights = 1/h_i)
```

```
## Dependent Var.:      lmwage      FGLSwage
##
## (Intercept)    -2.873*** (0.7390)  0.3995 (0.4141)
## educ           0.898*** (0.6813)  0.2965 (0.8316)
## exper          0.823*** (0.6121)  0.8319*** (0.0093)
## tenure         0.1693*** (0.0216)  0.1502*** (0.0244)
##
## S.E. type      IID      IID
## Observations   526      526
## R2              0.36842   0.21065
## Adj. R2        0.36244   0.20551
##
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

i. Realiza un test de White para los resultados del FGLS.

```
bptest(FGLSwage, varformulae ~ I(fitted(FGLSwage))^2, data=wage1)

##
## Studentized Breusch-Pagan test
##
## data: FGLSwage
## BP = 11.677, df = 2, p-value = 0.002913
```

Rechazamos hipótesis nula con 95% de confianza (p -value < 5). El test de White para FGLSwage concluye heterocedasticidad.

j. Estima FGLS con errores robustos de White-Huber. Reporta tus resultados y discute por que te gustaría estimar esto los errores estándar son muy diferentes a los obtenidos en h., por qué?

```
library(sandwich)
coefFGLS(FGLSwage, vcov = vcovHC(FGLSwage, types="HAC"))

##
## t test of coefficients:
##
## Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.898480  0.6169974  0.9668  0.344922
## educ        0.296580  0.8451588  0.4209  0.673618 ***
## exper       0.831864  0.0091421  3.4855  0.0005327 ***
## tenure      0.1502191  0.0279000  5.3842  1.182e-07 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Calcular FGLS con errores robustos permite ajustar los errores en dado caso de un mala especificación en la forma de la varianza, es decir, del vector de ponderación h_i . Si es así, $\frac{1}{h_i}$ aún es heterocedástico, por lo que el ajuste con errores robustos permite que los errores estándar sean válidos.

```
## FGLS robusto FGLS
## (Intercept)  0.6166 0.4141
## educ        0.8462 0.8316
## exper       0.8291 0.0093
## tenure      0.8273 0.0244
```

Los errores estándar no son muy diferentes. Lo anterior, ya que la ponderación h_i resolve parcialmente la heterocedasticidad y los errores robustos terminan de ajustarlos.

5. Ejecute el siguiente código de R

```
#Código
t <- c()
t.rob <- c()
# loop sampling and estimation
library(car)
for (i in 1:10000) {
  # sample data
  X <- 1:1000
  Y <- rnorm(n = 1000, mean = X, sd = 0.6 * X)

  # estimate regression model
  reg <- lm(Y ~ X)

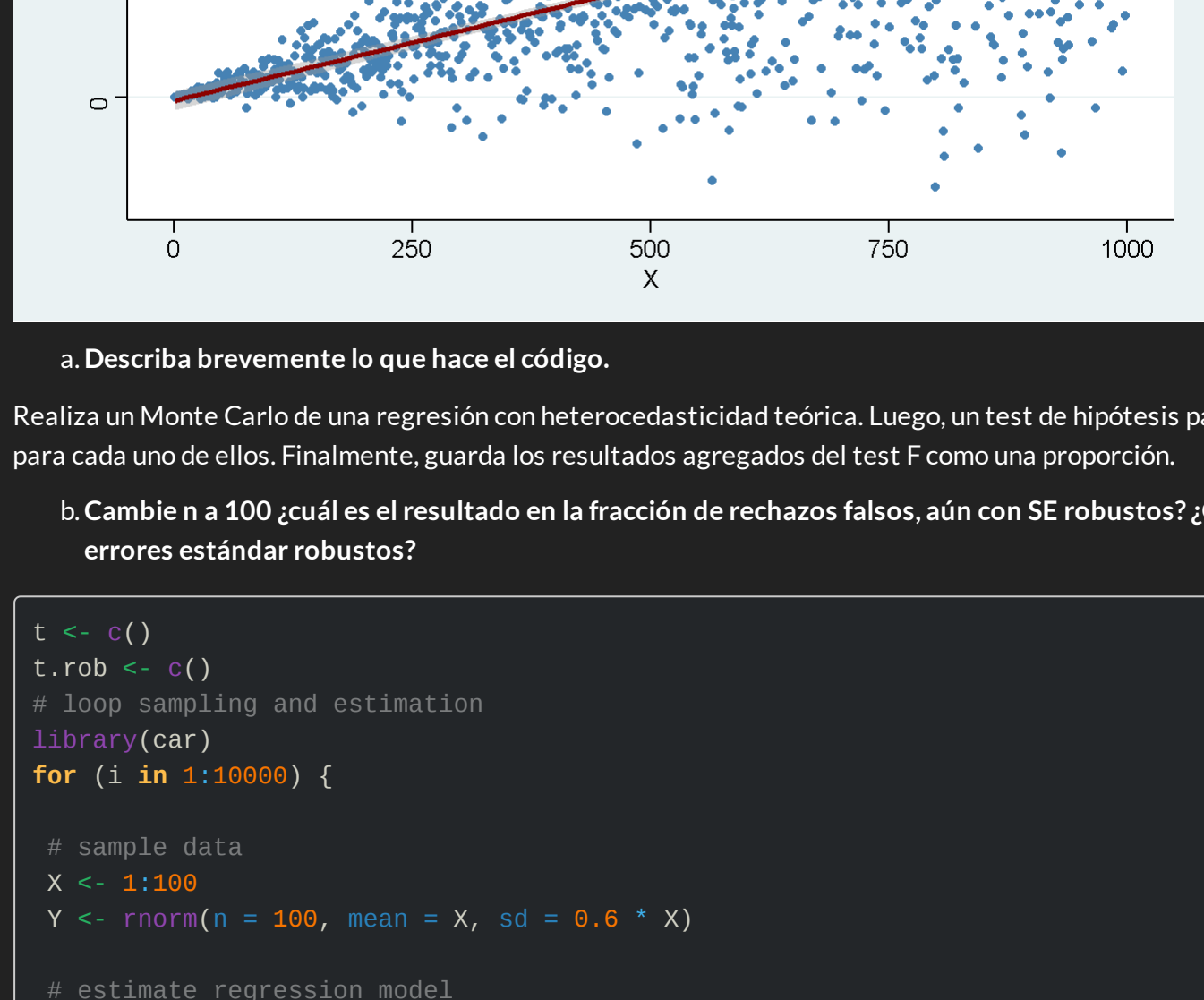
  # homoskedasticity-only significance test
  t[i] <- linearHypothesis(reg, "X = 1", white.adjust = "hcl")$Pr(>F) [2] < 0.05

  # robust significance test
  t.rob[i] <- linearHypothesis(reg, "X = 1", white.adjust = "hcl")$Pr(>F) [2] < 0.05
}
```

compute the fraction of false rejections
round(cbind(f = mean(t), t.rob = mean(t.rob)), 3)

```
##          t t.rob
## [1,] 0.972 0.951
```

```
#Hint
# generate heteroskedastic data
X <- 1:1000
Y <- rnorm(n = 1000, mean = X, sd = 0.6 * X)
# plot the data
ggplot(data=data.frame(X,Y), aes(X,Y)) + geom_point(color="steelblue") + geom_smooth(method=lm, color="darkred")
```



a. Describa brevemente lo que hace el código.

Realiza un Monte Carlo de una regresión con heterocedasticidad teórica. Luego, un test de hipótesis para β_1 y con errores robustos para cada uno de ellos. Finalmente, guarda los resultados agregados del test F como una proporción.

b. Cambia n a 100, ¿cuál es el resultado en la fracción de rechazos falsos, aún con SE robustos? ¿qué puede concluir sobre los errores estándar robustos?

```
t <- c()
t.rob <- c()
# loop sampling and estimation
library(car)
for (i in 1:10000) {
  # sample data
  X <- 1:100
  Y <- rnorm(n = 100, mean = X, sd = 0.6 * X)

  # estimate regression model
  reg <- lm(Y ~ X)

  # homoskedasticity-only significance test
  t[i] <- linearHypothesis(reg, "X = 1")$Pr(>F) [2] < 0.05

  # robust significance test
  t.rob[i] <- linearHypothesis(reg, "X = 1", white.adjust = "hcl")$Pr(>F) [2] < 0.05
}
```

compute the fraction of false rejections
round(cbind(f = mean(t), t.rob = mean(t.rob)), 3)

```
##          t t.rob
## [1,] 0.977 0.958
```

Con una muestra más pequeña, no rechaza la hipótesis nula en mayor proporción. Lo cual, dada la relación teórica, es un falso positivo. El cambio en la proporción es mayor para los errores robustos debido a sus propiedades asintóticas (aunque ambos menores que los no robustos). Es menos consistente.

c. Cambie SD a 1.6, ¿esto qué significa en términos de ajuste del modelo? ¿Qué efecto tiene en la fracción de rechazos falsos, aún con SE robustos?

```
t <- c()
t.rob <- c()
# loop sampling and estimation
library(car)
for (i in 1:10000) {
  # sample data
  X <- 1:1000
  Y <- rnorm(n = 1000, mean = X, sd = 1.6 * X)

  # estimate regression model
  reg <- lm(Y ~ X)

  # homoskedasticity-only significance test
  t[i] <- linearHypothesis(reg, "X = 1")$Pr(>F) [2] < 0.05

  # robust significance test
  t.rob[i] <- linearHypothesis(reg, "X = 1", white.adjust = "hcl")$Pr(>F) [2] < 0.05
}
```

```
##          t t.rob
## [1,] 0.977 0.961
```

Una mayor varianza aumenta la proporción de falsos positivos. El estimador es menos consistente. Aunque, en los tres casos, es conveniente no ignorar la heterocedasticidad y aplicar errores robustos, pues la proporción de falsos positivos será menor.

d. Describa formalmente el procedimiento para obtener WLS con estos datos.

$$u_i = 1.6x_i \rightarrow \text{Var} = \sigma u_i^2 = 1.6^2 x_i^2 = \sigma^2 x_i^2$$

Así,

$$h_i = x_i^2 \rightarrow w_i = 1/x_i^2$$

∴ el Mínimo Cuadrado Moderado es el siguiente:

$$\frac{\hat{y}_i}{x_i} = \beta_1 \frac{1}{x_i} + \beta_0 + \frac{u_i}{x_i}$$

Con $\frac{u_i}{x_i}$ homocedástico (= 1.6).

6. Estimamos un modelo de probabilidad lineal para arrestos de hombres jóvenes durante 1986.

$$\text{arr86} = \beta_0 + \beta_1 \cdot \text{pcnv} + \beta_2 \cdot \text{avgsen} + \beta_3 \cdot \text{tottime} + \beta_4 \cdot \text{ptime86} + \beta_5 \cdot \text{qemp86} + u$$

a. Usando los datos de `crime1`, estima este modelo por MCO y verifica que todos los valores ajustados están estrictamente entre 0 y 1. ¿Cuáles son los valores más chicos y grandes del estimado?

```
lmcrime <- lm(narr86 ~ pcnv + avgsen + tottime + ptime86 + qemp86, data=crime1)
sum(fitted(lmcrime) > 0) # 66 sum(fitted(lmcrime) < 1)
```

```
## [1] TRUE

min(fitted(lmcrime))

## [1] 0.88373395

max(fitted(lmcrime))

## [1] 0.9868572
```

b. Estima la ecuación por WLS como en la Sección 8-5.

$$\hat{h}_i = \hat{y}_i (1 - \hat{y}_i)$$

```
h.crime <- fitted.values(lmcrime) * (1 - fitted.values(lmcrime))
F.FGLScrime <- lm(narr86 ~ pcnv + avgsen + tottime + ptime86 + qemp86, data=crime1, weights= 1/h.crime)
```

c. Utiliza los estimados de FGLS para determinar si