

Lista Prática 4/5

Para os exercícios 1 a 3, usaremos a base de dados `RiceFarms` do pacote `splm` que tem informações sobre a produção de arroz na Indonésia. Você pode carregar a base usando:

```
1 # install.packages("splm")
2 data("RiceFarms", package = "splm")
```

É um painel com 1.026 observações, sendo $N = 171$ agricultores (`id`) em $T = 6$ períodos distintos (`time`). As variáveis relevantes são:

- `goutput`: output bruto de arroz em kg
- `seed`: sementes em kg
- `totlabor`: total de trabalhadores
- `size`: tamanho da plantação

Considere o seguinte modelo

$$\ln(\text{goutput}_{i,t}) = \alpha + \beta_1 \ln(\text{seed}_{i,t}) + \beta_2 \ln(\text{totlabor}_{i,t}) + \beta_3 \ln(\text{size}_{i,t}) + u_i + \nu_{i,t}$$

Exercício 1.

- Usando a função `plm()` do pacote `plm`, estime os modelos pelo estimador *between*, *pooled OLS*, *GLS* (efeitos aleatórios), e estimador *within* (efeitos fixos).
- Realize os testes de Breusch-Pagan/LM e de Hausman usando, respectivamente, as funções `plmtest()` e `phptest()`. Qual dos quatro estimadores utilizados no item (a) é mais adequado de acordo com os resultados dos testes?

Resposta:

```
a) library(dplyr)
2 library(plm)
3
4 # Carregando base de dados e transformando em pdata.frame
5 data("RiceFarms", package = "splm")
6 RF = pdata.frame(RiceFarms, index = c("id", "time"))
7
8 # Definindo modelos a serem estimados
9 models = c("within", "random", "pooling", "between")
10
11 # Estimando modelos via loop e atribuindo a objetos com nomes "variáveis" via
    assign()
12 for (x in models) {
13   assign(paste0("RF.", x), plm(log(goutput) ~ log(seed) + log(totlabor) + log(
    size), RF, model=x))
14 }
15
16 # Retornando as estimativas em tabela única
17 var_comuns = intersect(names(RF.within$coef), names(RF.random$coef))
18 rbind(within=RF.within$coef[var_comuns], random=RF.random$coef[var_comuns],
19       pooled=RF.pooling$coef[var_comuns], between=RF.between$coef[var_comuns])
```

```

1          log(seed) log(totlabor) log(size)
2 within  0.2095572      0.2891663 0.5023701
3 random   0.2199071      0.2855146 0.5278612
4 pooled   0.2220898      0.2842064 0.5333761
5 between  0.2280801      0.2766866 0.5506479

```

b) `# Realizando testes de Breusch-Pagan (Honda)`
`plmtest(RF.pooling, effect="individual")`

```

1 Lagrange Multiplier Test - (Honda) for balanced panels
2 data:  log(goutput) ~ log(seed) + log(totlabor) + log(size)
3 normal = 4.8396, p-value = 6.507e-07
4 alternative hypothesis: significant effects

```

Rejeitamos $H_0 : \sigma_u^2 = 0$ e, portanto, devemos considerar a presença de efeitos fixos
 \Rightarrow os modelos *between* e *pooled OLS* não são consistentes.

```

1 # Realizando teste de Hausman
2 phptest(RF.within, RF.random)

1 Hausman Test
2 data:  log(goutput) ~ log(seed) + log(totlabor) + log(size)
3 chisq = 3.775, df = 3, p-value = 0.2868
4 alternative hypothesis: one model is inconsistent

```

Não rejeitamos $H_0 : E(u|X) = 0$ (termo de erro individual é não-correlacionada com as covariadas) e, portanto, consideramos que ambos estimadores de *GLS* e *within* são consistentes \Rightarrow como *GLS* usa variações inter e intra-indivíduos para sua estimação, ele é mais eficiente e mais adequado do que *within* para este modelo/dados.

□

Exercício 2.

- a) *Estime analiticamente o modelo por estimador within (efeitos fixos).*
- b) *Estime analiticamente o modelo por GLS (efeitos aleatórios), calculando as variâncias dos termos de erro pelo método de Swamy e Arora (1972).*

Resposta:

- a) Note que a estimação analítica pode ser feita usando as fórmulas do *within* mostradas na seção 2.2.3 de Croissaint & Millo (2018) ou, também, transformando as variáveis via pré-multiplicação das variáveis pela matriz de transformação W e resolvendo por OLS. Aqui é importante prestar atenção nos graus de liberdade do estimador.

```

1 RiceFarms = RiceFarms %>% mutate(
2   # constant = 1, # within elimina o vetor de 1's
3   ln_goutput = log(goutput),
4   ln_seed = log(seed),
5   ln_totlabor = log(totlabor),
6   ln_size = log(size)
7 )
8

```

```

9 y = RiceFarms %>% select(ln_goutput) %>% as.matrix()
10 X = RiceFarms %>% select(ln_seed, ln_totlabor, ln_size) %>% as.matrix()
11 # Z = cbind(RiceFarms$constant, X) # within elimina o vetor de 1's
12
13 N = RiceFarms %>% select(id) %>% unique() %>% nrow()
14 T = RiceFarms %>% select(time) %>% unique() %>% nrow()
15 iota_T = rep(1, T)
16
17 # Calculando matrizes de transformação B e W
18 B = diag(N) %x% (iota_T %*% solve(t(iota_T) %*% iota_T) %*% t(iota_T))
19 W = diag(N*T) - B
20
21 # vetor de estimativas gamma_hat = (alpha, beta)
22 gamma_hat = solve(t(X) %*% W %*% X) %*% t(X) %*% W %*% y
23
24 # valores ajustados e erros
25 y_hat = X %*% gamma_hat
26 e_hat = y - y_hat
27
28 ## Estimando variancia do termo de erro
29 sigma2_v = t(e_hat) %*% W %*% e_hat / (N*T - ncol(X) - N) # NT - K - N g.l.
30
31 ## Estimando a matriz de variancia/covariancia das estimativas gamma
32 vcov_hat = c(sigma2_v) * solve(t(X) %*% W %*% X)
33
34 ## Calculando erros padrao das estimativas gamma
35 std_error = sqrt(diag(vcov_hat)) # Raiz da diagonal da matriz de covariâncias
36
37 ## Calculando estatisticas t das estimativas gamma
38 t_stat = gamma_hat / std_error
39
40 ## Calculando p-valores das estimativas gamma
41 p_value = 2 * pt(q = -abs(t_stat), df = N*T - ncol(X) - N) # NT - K - N g.l.
42
43 ## Organizando os resultados da regressao em uma matriz
44 results = cbind(gamma_hat, std_error, t_stat, p_value)
45
46 ## Nomeando as colunas da matriz de resultados
47 colnames(results) = c("Estimate", "Std. Error", "t stat", "Pr(>|t|)")
48 results

```

```

1      Estimate Std. Error    t stat    Pr(>|t|)
2 ln_seed      0.2095572  0.03182723  6.584212 7.989993e-11
3 ln_totlabor  0.2891663  0.03549156  8.147465 1.319048e-15
4 ln_size      0.5023701  0.03801667 13.214469 2.114140e-36

```

Comparando com os resultados obtidos no Exercício 1:

```

1 summary(RF.within)$coef # comparando com estimado via plm()

```

```

1      Estimate Std. Error    t-value    Pr(>|t|)
2 log(seed)      0.2095572  0.03182723  6.584212 7.989993e-11
3 log(totlabor)  0.2891663  0.03549156  8.147465 1.319048e-15
4 log(size)      0.5023701  0.03801667 13.214469 2.114140e-36

```

- b) Lembre-se que, na estimação GLS (efeitos aleatórios), usamos resíduos de outro estimador para calcular $\hat{\sigma}_v$, $\hat{\sigma}_u$ e, conseqüentemente, $\hat{\Omega}$.

Para usar o método de Swamy e Arora (1972) precisamos obter os erros das estimações *within* e *between* para calcular:

$$\hat{\sigma}_l^2 = \frac{\hat{\varepsilon}'_B B \hat{\varepsilon}_B}{N - K - 1}, \quad \hat{\sigma}_v^2 = \frac{\hat{\varepsilon}'_W W \hat{\varepsilon}_W}{N(T - 1) - K} \quad \text{e} \quad \hat{\sigma}_u^2 = \frac{\hat{\sigma}_l^2 - \hat{\sigma}_v^2}{T}$$

```

1 # obtendo estimativas within e between
2 gamma_hat_w = solve(t(X) %*% W %*% X) %*% t(X) %*% W %*% y
3 gamma_hat_b = solve(t(Z) %*% B %*% Z) %*% t(Z) %*% B %*% y
4
5 # obtendo os valores ajustados e resíduos within e between
6 y_hat_w = X %*% gamma_hat_w
7 e_w = y - y_hat_w
8
9 y_hat_b = Z %*% gamma_hat_b
10 e_b = y - y_hat_b
11
12 # Calculando os termos de erro
13 sigma2_l = (t(e_b) %*% B %*% e_b) / (N - ncol(X) - 1)
14 sigma2_v = (t(e_w) %*% W %*% e_w) / (N * (T-1) - ncol(X))
15 sigma2_u = (sigma2_l + sigma2_v) / T
16
17 sigmas2 = cbind(sigma2_l, sigma2_v, sigma2_u)
18 colnames(sigmas2) = c("sigma2_l", "sigma2_v", "sigma2_u")
19 sigmas2

```

```

1 sigma2_l sigma2_v sigma2_u
2 0.2204066 0.1323742 0.0587968

```

Agora, calculamos Ω^{-1} e para obter

$$\hat{\gamma}_{GLS} = (Z'\Omega^{-1}Z)^{-1}(Z'\Omega^{-1}y) = \quad \text{e} \quad V(\hat{\gamma}_{GLS}) = (Z'\Omega^{-1}Z)^{-1}$$

```

1 # Calculando \Omega^{-1}
2 Omega_1 = c(sigma2_l^(-1)) * B + c(sigma2_v^(-1)) * W
3 dim(Omega_1) # NT x NT
4
5 # vetor de estimativas gamma_hat = (alpha, beta)
6 gamma_hat = solve(t(Z) %*% Omega_1 %*% Z) %*% (t(Z) %*% Omega_1 %*% y)
7 gamma_hat
8
9 ## Estimando a matriz de variancia/covariancia das estimativas gamma
10 vcov_hat = solve(t(Z) %*% Omega_1 %*% Z)
11
12 ## Calculando erros padrao das estimativas gamma
13 std_err = sqrt(diag(vcov_hat)) # Raiz da diagonal da matriz de covariâncias
14
15 ## Calculando estatisticas t das estimativas gamma
16 t_stat = gamma_hat / std_err
17
18 ## Calculando p-valores das estimativas gamma
19 p_value = 2 * pt(q = -abs(t_stat), df = nrow(Z) - ncol(Z)) # NT - K - 1 g.l.
20
21 ## Organizando os resultados da regressao em uma matriz
22 results_swar = cbind(gamma_hat, std_err, t_stat, p_value)
23

```

```

24 ## Nomeando as colunas da matriz de resultados
25 colnames(results_swar) = c("Estimate", "Std. Error", "t stat", "Pr(>|t|)")
26 rownames(results_swar) = c("(Intercept)", "ln_seed", "ln_totlabor", "ln_size")
27 results_swar

```

```

1      Estimate Std. Error    t stat    Pr(>|t|)
2 (Intercept)  5.3123104  0.20415280  26.021247  6.068667e-115
3 ln_seed      0.2199071  0.02837609   7.749733  2.214702e-14
4 ln_totlabor  0.2855146  0.03110824   9.178101  2.366789e-19
5 ln_size      0.5278612  0.03268695  16.148989  2.010632e-52

```

Comparando com os resultados obtidos no Exercício 1:

```

1 summary(RF.random)$coef # comparando com estimado via plm()

```

```

1      Estimate Std. Error    z-value    Pr(>|z|)
2 (Intercept)  5.3123104  0.20422998  26.011413  3.678718e-149
3 log(seed)    0.2199071  0.02838682   7.746804  9.423402e-15
4 log(totlabor) 0.2855146  0.03112000   9.174632  4.531184e-20
5 log(size)    0.5278612  0.03269931  16.142887  1.274387e-58

```

□

Exercício 3. *Estime numericamente o modelo por máxima verossimilhança (FIML) usando os seguintes passos:*

- (1) Chutar valores iniciais para as estimativas $\hat{\gamma}$ (pode usar tudo igual a 0)
- (2) Mostrar a iteração e as estimativas atuais de α e β 's
- (3) Obter $\hat{\varepsilon} = y - \hat{y}$ e calcular $\hat{\phi}^2$ usando

$$\hat{\phi}^2 = \frac{\hat{\varepsilon}'W\hat{\varepsilon}}{(T-1)\hat{\varepsilon}'B\hat{\varepsilon}} \quad (3.14)$$

- (4) Calcular $\hat{\sigma}_\nu^2$ usando

$$\hat{\sigma}_\nu^2 = \frac{\hat{\varepsilon}'W\hat{\varepsilon} + \hat{\phi}^2\hat{\varepsilon}'B\hat{\varepsilon}}{NT} \quad (3.13)$$

- (5) Obter $\hat{\Omega}^{-1/2}$ a partir da relação

$$\hat{\Omega}^{-1/2} = \frac{1}{\hat{\sigma}_l}B + \frac{1}{\hat{\sigma}_\nu}W \iff \hat{\Omega}^{-1/2} = \frac{\hat{\phi}B + W}{\hat{\sigma}_\nu} \quad (2.29')$$

em que $\hat{\phi} = \hat{\sigma}_\nu/\hat{\sigma}_l$.

- (6) Calcular o novas estimativas $\hat{\gamma}'$ usando

$$\hat{\gamma}' = (\tilde{Z}'\tilde{Z})^{-1}\tilde{Z}'\tilde{y} \quad (3.12)$$

tal que $\tilde{Z} = \hat{\Omega}^{-1/2}Z$, e $\tilde{y} = \hat{\Omega}^{-1/2}y$

(7) Verificar convergência das estimativas de acordo com:

$$\text{distância} = \max\{abs(\hat{\gamma}' - \hat{\gamma})\} < 1 \times 10^{-10} = \text{tolerância}$$

Se não convergiu (i.e., expressão acima não foi satisfeita), volte ao passo (2), definindo $\hat{\gamma} \equiv \hat{\gamma}'$ e iniciando uma nova iteração

(8) Obter, mais uma vez, $\hat{\phi}^2$ e $\hat{\sigma}_\nu^2$ para calcular

$$\hat{\sigma}_l = \frac{\hat{\sigma}_\nu}{\hat{\phi}} \quad e \quad \hat{\sigma}_u = \sqrt{\frac{\hat{\sigma}_l^2 - \hat{\sigma}_\nu^2}{T}}$$

(9) Calcular $V(\hat{\gamma})$ usando¹

$$V(\hat{\gamma}) = \left(\frac{1}{\hat{\sigma}_\nu^2} Z' W Z + \frac{1}{\hat{\sigma}_l^2} Z' B Z \right)^{-1}$$

(10) Obter erros padrão das estimativas, estatísticas t e p -valores. Não é necessário fazer esse último passo para $\hat{\sigma}_\nu$ e $\hat{\sigma}_u$

(11) Comparar os resultados obtidos com a estimação via função `pglm()` do pacote `pglm`²

Resposta:

```
1 # Iremos realizar iterações até a convergência de \hat{\gamma}
2 # (1) Chutar valores iniciais para as estimativas \hat{\gamma}
3 gamma_ini = c(0, 0, 0, 0) # chute inicial
4
5 tol = 1e-10 # tolerância para convergência
6 dist = 1 # distância inicial - apenas para entrar no loop/while
7 it = 0 # número de iterações
8
9 while (dist > tol) {
10   # (2) Mostrar a iteração atual e as estimativas atuais de \alpha e \beta's
11   print(paste0("iteração ", it,
12     ": alpha = ", round(gamma_ini[1], 5),
13     " | b1 = ", round(gamma_ini[2], 5),
14     " | b2 = ", round(gamma_ini[3], 5),
15     " | b3 = ", round(gamma_ini[4], 5)
16 ))
17
18   # (3) Obter \hat{\varepsilon} = y - \hat{y} e calcular \hat{\phi}^2
19   y_hat = Z %% gamma_ini
20   e = as.vector(y - y_hat)
21
22   phi2_hat = (t(e) %% W %% e) / ((T-1) * (t(e) %% B %% e))
23   phi_hat = as.vector(sqrt(phi2_hat))
24
25   # (4) Calcular \hat{\sigma}^2_v
26   sigma2_v = (t(e) %% W %% e + phi2_hat * t(e) %% B %% e) / (N*T)
```

¹Caso retorne “Error in (...): non-conformable arrays”, talvez seja necessário transformar os números absolutos em vetores (via função `as.vector()` ou `c()`).

²Não ficará exatamente igual, mas os valores são bem próximos.

```

27 sigma_v = as.vector(sqrt(sigma2_v))
28
29 # (5) Obter  $\hat{\Omega}^{-1/2}$ 
30 Omega_sqrt = (phi_hat * B + W) / sigma_v
31
32 # (6) Calcular o novas estimativas  $\hat{\gamma}$ 
33 Z_til = Omega_sqrt %>% Z
34 y_til = Omega_sqrt %>% y
35
36 gamma_fim = solve(t(Z_til) %>% Z_til) %>% t(Z_til) %>% y_til
37
38 # (7) Verificar convergência das estimativas
39 dist = max(abs(gamma_fim - gamma_ini)) # calculando distância
40 gamma_ini = gamma_fim
41 it = it + 1
42 }

1 [1] "iteração 0: alpha = 0 | b1 = 0 | b2 = 0 | b3 = 0"
2 [1] "iteração 1: alpha = 5.28351 | b1 = 0.20962 | b2 = 0.28915 | b3 = 0.50252"
3 [1] "iteração 2: alpha = 5.31177 | b1 = 0.21976 | b2 = 0.28559 | b3 = 0.5275"
4 [1] "iteração 3: alpha = 5.31253 | b1 = 0.21996 | b2 = 0.28548 | b3 = 0.528"
5 [1] "iteração 4: alpha = 5.31254 | b1 = 0.21997 | b2 = 0.28548 | b3 = 0.52801"
6 [1] "iteração 5: alpha = 5.31254 | b1 = 0.21997 | b2 = 0.28548 | b3 = 0.52801"
7 [1] "iteração 6: alpha = 5.31254 | b1 = 0.21997 | b2 = 0.28548 | b3 = 0.52801"

1 # (8) Calcular  $\hat{\phi}^2$ ,  $\hat{\sigma}_v^2$ ,  $\hat{\sigma}_l$  e  $\hat{\sigma}_u$ 
2 y_hat = Z %>% gamma_ini
3 e = as.vector(y - y_hat)
4 phi2_hat = (t(e) %>% W %>% e) / ((T-1) * (t(e) %>% B %>% e))
5 phi_hat = sqrt(phi2_hat)
6
7 sigma2_v = (t(e) %>% W %>% e + phi2_hat * t(e) %>% B %>% e) / (N*T)
8 print(paste0("sd_idios = ", round(sqrt(sigma2_v), 6)))
9
10 sigma_l = sqrt(sigma2_v) / phi_hat
11 sigma2_u = (sigma_l^2 - sigma2_v) / T
12
13 # (9) Calcular  $V(\hat{\gamma})$ 
14 V = solve(c(1/sigma_v^2) * t(Z) %>% W %>% Z + c(1/sigma_l^2) * t(Z) %>% B %>% Z)
15 std_error = sqrt(diag(V))
16 t_value = gamma_ini / std_error
17 p_value = pt(-abs(t_value), df = N*T - ncol(Z)) # NT - K - 1
18
19 # (10) Resultados
20 result = matrix(NA, nrow=(ncol(Z) + 2), ncol=4)
21 colnames(result) = c("Estimate", "Std. Error", "t stat", "Pr(>|t|)")
22 rownames(result) = c("(Intercept)", "ln_seed", "ln_totlabor", "ln_size",
23 "sd.id", "sd.idios")
24 result[1:ncol(Z),] = cbind(gamma_ini, std_error, t_value, p_value)
25 result[(ncol(Z)+1):(ncol(Z)+2),1] = rbind(sqrt(sigma2_u), sigma_v)
26 result
27
28 # (11) Observe que obtivemos resultados próximos do obtido via 'pglm()'.
29 pglm::pglm(log(goutput) ~ log(seed) + log(totlabor) + log(size),
30 RiceFarms, family = "gaussian") %>% summary() %>% coef()

1 Estimate Std. Error t stat Pr(>|t|)
2 (Intercept) 5.3125397 0.20376061 26.072457 1.359514e-115
3 ln_seed 0.2199672 0.02832543 7.765714 9.831510e-15

```

```

4 ln_totlabor 0.2854829 0.03104536 9.195671 1.017714e-19
5 ln_size     0.5280116 0.03262143 16.186035 6.228455e-53
6 sd.id       0.1190403          NA          NA          NA
7 sd.idios    0.3636631          NA          NA          NA

1 # Comparando os resultados com os obtidos por pglm()
2 pglm::pglm(log(goutput) ~ log(seed) + log(totlabor) + log(size),
3            RiceFarms, family = "gaussian") %>% summary() %>% coef()

1              Estimate Std. error  t value      Pr(> t)
2 (Intercept)  5.3125396 0.203770564 26.071183 7.739884e-150
3 log(seed)    0.2199672 0.028330404 7.764349 8.206540e-15
4 log(totlabor) 0.2854829 0.031046602 9.195304 3.739365e-20
5 log(size)    0.5280116 0.032648677 16.172526 7.880042e-59
6 sd.id        0.1190406 0.017129023 6.949645 3.662067e-12
7 sd.idios     0.3636625 0.008600951 42.281665 0.000000e+00

```

□

Para os exercícios 4 e 5, usaremos a base de dados `card` do pacote `wooldridge` que foi usado por David Card (1995). Você pode carregar a base usando:

```

1 # install.packages("wooldridge")
2 data("card", package="wooldridge")

```

Possui 3.010 observações e as variáveis relevantes são:

- `lwage`: log do salário
- `educ`: anos de educação
- `exper`: experiência (idade - anos de educação - 6)
- `expersq`: experiência²
- `black`: = 1 se pessoa negra
- `nearc4`: = 1 se próximo de faculdade (4-year college)
- `fatheduc`: anos de educação do pai
- `motheduc`: anos de educação da mãe

Exclua as observações com valores ausentes (NA's) usando `na.omit()` e considere o modelo:

$$\text{lwage}_i = \alpha + \beta_1 \text{educ}_i + \beta_2 \text{exper}_i + \beta_3 \text{expersq}_i + \beta_4 \text{black}_i + \varepsilon_i$$

Exercício 4. Considere que a variável `nearc4` é o instrumento de `educ`.

- Estime o modelo por OLS usando a função `lm()`.
- Estime o modelo por IV usando a função `ivreg()`³ do pacote `AER`. Resuma os resultados usando `summary(..., diagnostics=TRUE)` para fazer os testes de instrumentos fracos, Wu-Hausman e Sargan. Analise-os.

³A sintaxe da função é `ivreg(y ~ X + W | Z + W, data)`, em que X representa uma variável endógena, Z é o seu instrumento, e W é uma variável exógena. Note que a inserção da fórmula é parecida com da função `lm()`, porém inclui-se uma barra vertical para inserir as variáveis instrumentais (os instrumentos das variáveis exógenas são elas mesmas).

c) Estime o modelo por 2SLS usando a função `lm()`

d) Estime o modelo por IV de forma analítica. Observe que as matrizes

$$X_{N \times (K+1)} = \begin{bmatrix} 1 & x_{11}^* & x_{12} & \cdots & x_{1K} \\ 1 & x_{21}^* & x_{22} & \cdots & x_{2K} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1}^* & x_{N2} & \cdots & x_{NK} \end{bmatrix} \quad e \quad Z_{N \times (K+1)} = \begin{bmatrix} 1 & z_{11}^* & x_{12} & \cdots & x_{1K} \\ 1 & z_{21}^* & x_{22} & \cdots & x_{2K} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & z_{N1}^* & x_{N2} & \cdots & x_{NK} \end{bmatrix}$$

possuem a mesma dimensão, e a variável endógena $x_{\cdot 1}^*$ e o seu instrumento $z_{\cdot 1}^*$ estão na mesma posição (coluna) em suas respectivas matrizes. Lembre-se também que:

$$\hat{\beta}_{IV} = (Z'X)^{-1}Z'y \quad e \quad V(\hat{\beta}) = \hat{\sigma}^2(X'P_ZX)^{-1}$$

Resposta:

```
a) # David Card (1995), Using Geographic Variation in College Proximity to
# Estimate the Return to Schooling, in Aspects of Labour Market Behavior
library(AER)
library(dplyr)
data("card", package="wooldridge")
card = na.omit(card) # Omitindo NA's

# OLS
ols = lm(lwage ~ educ + exper + expersq + black, data=card)
summary(ols)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.6644195	0.1040757	44.818	< 2e-16 ***
educ	0.0806292	0.0052549	15.344	< 2e-16 ***
exper	0.0960991	0.0111488	8.620	< 2e-16 ***
expersq	-0.0028008	0.0005971	-4.691	2.95e-06 ***
black	-0.1610134	0.0321305	-5.011	6.01e-07 ***

```
b) # IV
iv = ivreg(lwage ~ educ + exper + expersq + black
| nearc4 + exper + expersq + black, data=card)
summary(iv, diagnostics=TRUE)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.996941	1.208832	1.652	0.098740 .
educ	0.229317	0.067314	3.407	0.000674 ***
exper	0.184291	0.042026	4.385	1.24e-05 ***
expersq	-0.004757	0.001146	-4.152	3.47e-05 ***
black	-0.068451	0.057363	-1.193	0.232935

Diagnostic tests:

	df1	df2	statistic	p-value
Weak instruments	1	1595	14.734	0.000129 ***
Wu-Hausman	1	1594	7.426	0.006500 **
Sargan	0	NA	NA	NA

- *Instrumentos fracos*: é um teste F nos instrumentos no 1º estágio, em que a hipótese nula é que os instrumentos são fracos e, portanto, rejeitá-la significa que os instrumentos não são fracos.
- *Wu-Hausman*: Teste a consistência das estimativas OLS, sob a hipótese que a estimação por IV é consistente. Quando rejeitamos, significa que OLS não é consistente, sugerindo presença de endogeneidade. E quando aceitamos a nula, significa que ambas estimativas são similares, e a endogeneidade pode não ser um grande problema.
- *Sargan (teste J)*: É um teste de exogeneidade dos instrumentos usando restrições sobre-identificadas. Só pode ser utilizado quando houver mais instrumentos do que regressores endógenos. Se a nula é rejeitada, significa que pelo menos um dos instrumentos é inválido.

```
c) # 2SLS
2 stage1 = lm(educ ~ nearc4 + exper + expersq + black, data=card)
3 educ_hat = stage1$fitted.values
4
5 stage2 = lm(lwage ~ educ_hat + exper + expersq + black, data=card)
6 summary(stage2)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.9969409	1.0516122	1.899	0.0578 .
educ_hat	0.2293167	0.0585591	3.916	9.38e-05 ***
exper	0.1842915	0.0365605	5.041	5.17e-07 ***
expersq	-0.0047573	0.0009968	-4.773	1.98e-06 ***
black	-0.0684507	0.0499026	-1.372	0.1704

```
d) # Criando matrizes
2 X = card %>% mutate(constant = 1) %>%
3   select(constant, educ, exper, expersq, black) %>% as.matrix()
4 Z = card %>% mutate(constant = 1) %>%
5   select(constant, nearc4, exper, expersq, black) %>% as.matrix()
6 y = card$lwage %>% as.matrix()
7
8 head(X, 4)
9 head(Z, 4)
```

```
1 constant   educ  exper expersq black
2           1    12     9      81    0
3           1    12    16     256    0
4           1    11    10     100    0
5           1    12    16     256    0
6
7 constant nearc4  exper expersq black
8           1     0     9      81    0
9           1     0    16     256    0
10          1     1    10     100    0
11          1     1    16     256    0
```

```
1 # Estimando o modelo IV analiticamente
2 beta_hat = solve(t(Z) %*% X) %*% t(Z) %*% y
3
4 y_hat = X %*% beta_hat
5 e = y - y_hat
6
7 sigma2_hat = t(e) %*% e / (nrow(X) - ncol(X))
8 vcov_hat = c(sigma2_hat)*solve(t(X) %*% Z %*% solve(t(Z) %*% Z) %*% t(Z) %*% X)
```

```

9
10 std_error = sqrt(diag(vcov_hat))
11 t_stat = beta_hat / std_error
12 p_value = 2 * pt(q = -abs(t_stat), df = nrow(X) - ncol(X))
13
14 results = cbind(beta_hat, std_error, t_stat, p_value)
15 colnames(results) = c("Estimate", "Std. Error", "t stat", "Pr(>|t|)")
16 results %>% round(6)

```

	Estimate	Std. Error	t stat	Pr(> t)
1 constant	1.996941	1.208832	1.651959	0.098740
2 educ	0.229317	0.067314	3.406675	0.000674
4 exper	0.184291	0.042026	4.385131	0.000012
5 expersq	-0.004757	0.001146	-4.152044	0.000035
6 black	-0.068451	0.057363	-1.193285	0.232935

□

Exercício 5. Agora, considere que *educ* tem três instrumentos: *nearc4*, *fatheduc* e *motheduc*.

- Estime o modelo por IV usando a função *ivreg()* e analise os testes de instrumentos fracos, Wu-Hausman e Sargan.
- Estime o modelo por IV de forma analítica. Dica: consulte as notas de aula sobre a forma de estimação IV com modelo sobre-identificado.

Resposta:

```

a) # IV
2 iv2 = ivreg(lwage ~ educ + exper + expersq + black
3           | nearc4 + fatheduc + motheduc + exper + expersq + black,
4           data=card)
5 summary(iv2, diagnostics=TRUE)

```

	Estimate	Std. Error	t value	Pr(> t)
1 (Intercept)	3.9226903	0.3175946	12.351	< 2e-16 ***
2 educ	0.1219738	0.0175248	6.960	4.95e-12 ***
4 exper	0.1206222	0.0150691	8.005	2.29e-15 ***
5 expersq	-0.0033448	0.0006469	-5.170	2.63e-07 ***
6 black	-0.1352751	0.0343561	-3.937	8.59e-05 ***

7

8 Diagnostic tests:

	df1	df2	statistic	p-value
10 Weak instruments	3	1593	54.705	<2e-16 ***
11 Wu-Hausman	1	1594	6.399	0.0115 *
12 Sargan	2	NA	5.864	0.0533 .

- Havendo um modelo sobre-identificado (n^0 de instrumentos $> n^0$ de variáveis endógenas), podemos criar uma nova variável instrumental por meio da combinação linear dos instrumentos para usá-la como única variável instrumental da variável endógena $x_{.1}^*$. Para isto, podemos projetar as variáveis instrumentais na variável endógena para obter os “pesos” de cada instrumento na nova variável, a partir do modelo:

$$x_{.1}^* = \gamma_0 + \gamma_1 z_{.1a}^* + \gamma_2 z_{.1b}^* + \gamma_3 z_{.1c}^* + \varepsilon$$

Agora, usamos as estimativas $\hat{\gamma}$ para criar o novo instrumento:

$$z_{.1}^* = \hat{x}_{.1}^* = \hat{\gamma}_0 + \hat{\gamma}_1 z_{.1a}^* + \hat{\gamma}_2 z_{.1b}^* + \hat{\gamma}_3 z_{.1c}^*$$

```
1 ## Estimação IV Analítica
2 # Criando nova variável instrumental
3 x_star = card$educ %>% as.matrix()
4 Z_star = card %>% mutate(constant = 1) %>%
5   select(constant, nearc4, fatheduc, motheduc) %>%
6   as.matrix()
7
8 gamma = solve(t(Z_star) %*% Z_star) %*% t(Z_star) %*% x_star
9 gamma
10
11 z_star = Z_star %*% gamma # combinação linear dos instrumentos
```

```
1      [,1]
2 constant 10.2255758
3 nearc4    0.3429114
4 fatheduc  0.1888914
5 motheduc  0.1500375
```

Agora, os passos são praticamente os mesmos realizados no Exercício 4, tendo como diferença a inserção da nova variável instrumental na matrix Z .

```
1 # Criando matrizes
2 X = card %>% mutate(constant = 1) %>%
3   select(constant, educ, exper, expersq, black) %>% as.matrix()
4 Z = cbind(z_star, card) %>% mutate(constant = 1) %>%
5   select(constant, z_star, exper, expersq, black) %>% as.matrix()
6 y = card$lwage %>% as.matrix()
7
8 head(X, 4)
9 head(Z, 4)
```

```
1  constant    educ  exper  expersq  black
2      1      12      9      81      0
3      1      12     16     256      0
4      1      11     10     100      0
5      1      12     16     256      0
6
7  constant z_star  exper  expersq  black
8      1 12.937      9      81      0
9      1 14.670     16     256      0
10     1 14.446     10     100      0
11     1 13.129     16     256      0
```

```
1 # Estimando o modelo IV sobre-identificado
2 beta_hat = solve(t(Z) %*% X) %*% t(Z) %*% y
3
4 y_hat = X %*% beta_hat
5 e = y - y_hat
6
7 sigma2_hat = t(e) %*% e / (nrow(X) - ncol(X))
8 vcov_hat = c(sigma2_hat)*solve(t(X) %*% Z %*% solve(t(Z) %*% Z) %*% t(Z) %*% X)
9
10 std_error = sqrt(diag(vcov_hat))
11 t_stat = beta_hat / std_error
```

```

12 p_value = 2 * pt(q = -abs(t_stat), df = nrow(X) - ncol(X))
13
14 results = cbind(beta_hat, std_error, t_stat, p_value)
15 colnames(results) = c("Estimate", "Std. Error", "t_stat", "Pr(>|t|)")
16 results

```

	Estimate	Std. Error	t_stat	Pr(> t)
constant	3.982329	0.317704	12.534707	0.0e+00
educ	0.118649	0.017532	6.767579	0.0e+00
exper	0.118650	0.015049	7.884154	0.0e+00
expersq	-0.003301	0.000645	-5.115249	0.0e+00
black	-0.137345	0.034268	-4.007911	6.4e-05

□