

Empirical Exercise - E6.1

Chi-Yuan Fang

2021-03-21

Use the **Birthweight_Smoking** data set introduced in Empirical Exercise E5.3 to answer the following questions.

- a. Regress *Birthweight* on *Smoker*. What is the estimated effect of smoking on birth weight?

Solution

```
library(ggplot2); library(dplyr);
library(jtools); library(ggstance);

##
## Attaching package: 'ggstance'
## The following objects are masked from 'package:ggplot2':
##
##   geom_errorbarh, GeomErrorbarh
library(broom.mixed); library(huxtable)

## Warning in checkMatrixPackageVersion(): Package version inconsistency detected.
## TMB was built with Matrix version 1.3.2
## Current Matrix version is 1.2.18
## Please re-install 'TMB' from source using install.packages('TMB', type = 'source') or ask CRAN for a
## Registered S3 method overwritten by 'broom.mixed':
##   method      from
##   tidy.gamlss broom
##
## Attaching package: 'huxtable'
## The following object is masked from 'package:dplyr':
##
##   add_rownames
## The following object is masked from 'package:ggplot2':
##
##   theme_grey
# import data
library(readxl)
Birthweight_Smoking <- read_excel("Birthweight_Smoking/Birthweight_Smoking.xlsx")

E61a <- lm(birthweight ~ smoker, data = Birthweight_Smoking)

summ(E61a, confint = TRUE, digits = 4)

## MODEL INFO:
## Observations: 3000
```

```
## Dependent Variable: birthweight
## Type: OLS linear regression
##
## MODEL FIT:
## F(1,2998) = 88.2793, p = 0.0000
## R2 = 0.0286
## Adj. R2 = 0.0283
##
## Standard errors: OLS
## -----
##               Est.      2.5%      97.5%      t val.      p
## -----
## (Intercept)    3432.0600   3408.7840   3455.3359   289.1154   0.0000
## smoker         -253.2284   -306.0736   -200.3831   -9.3957   0.0000
## -----
```

- b. Regress *Birthweight* on *Smoker*, *Alcohol*, and *Nprevist*.
- Using the two conditions in Key Concept 6.1, explain why the exclusion of *Alcohol* and *Nprevist* could lead to omitted variable bias in the regression estimated in (a).
 - Is the estimated effect of smoking on birth weight substantially different from the regression that excludes *Alcohol* and *Nprevist*? Does the regression in (a) seem to suffer from omitted variable bias?
 - Jane smoked during her pregnancy, did not drink alcohol, and had 8 prenatal care visits. Use the regression to predict the birth weight of Jane's child.
 - Compute R^2 and \bar{R}^2 . Why are they so similar?
 - How should you interpret the coefficient on *Nprevist*? Does the coefficient measure a causal effect of prenatal visits on birth weight? If not, what does it measure?

Solution

```
# nprevist: total number of prenatal visits
E61b <- lm(birthweight ~ smoker + alcohol + nprevist, data = Birthweight_Smoking)

summ(E61b, confint = TRUE, digits = 4)
```

```
## MODEL INFO:
## Observations: 3000
## Dependent Variable: birthweight
## Type: OLS linear regression
##
## MODEL FIT:
## F(3,2996) = 78.4697, p = 0.0000
## R2 = 0.0729
## Adj. R2 = 0.0719
##
## Standard errors: OLS
## -----
##               Est.      2.5%      97.5%      t val.      p
## -----
## (Intercept)    3051.2486   2984.5516   3117.9456   89.7005   0.0000
## smoker         -217.5801   -269.8923   -165.2679   -8.1553   0.0000
## alcohol        -30.4913   -179.9677    118.9851   -0.4000   0.6892
## nprevist        34.0699    28.4720    39.6679   11.9334   0.0000
## -----
```

i.

- Smoking may be correlated with both alcohol and the number of pre-natal doctor visits.
 - Both alcohol consumption and the number of doctor visits may have their own independent affects on birthweight.
- ii. Because the estimated is smaller than (a), the regression in (a) may suffer from omitted variable bias.

```
export_summs(E61a, E61b,
             model.names = c("Model (a)", "Model (b)"))
```

	Model (a)	Model (b)
(Intercept)	3432.06 *** (11.87)	3051.25 *** (34.02)
smoker	-253.23 *** (26.95)	-217.58 *** (26.68)
alcohol		-30.49 (76.23)
nprevist		34.07 *** (2.85)
N	3000	3000
R2	0.03	0.07

*** p < 0.001; ** p < 0.01; * p < 0.05.

iii.

```
# predict value
E61biii <- function(x){
  E61b$coefficients %*% matrix(c(1, x), ncol = 1)
}
```

```
# predict value:
# smoker = 1, alcohol = 0, nprevist = 8
E61biii(c(1, 0, 8))
```

```
##           [,1]
## [1,] 3106.228
```

- iv. Because the sample size is very large, they are nearly identical.
- v. *Nprevist* is a control variable. It captures mother's access to healthcare and health. Thus, its coefficient does not have a causal interpretation.
- c. Estimate the coefficient on Smoking for the multiple regression model in (b), using the three-step process in Appendix 6.3 (the Frisch–Waugh theorem). Verify that the three-step process yields the same estimated coefficient for *Smoking* as that obtained in (b).

Solution

The OLS estimator in multiple regression can be computed by a sequence of shorter regressions. Consider the multiple regression model

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i, \quad i = 1, \dots, n. \quad (1)$$

The OLS estimator of β_1 can be computed in three steps:

1. Regress X_1 on X_2, X_3, \dots, X_k , and let \tilde{X}_1 denote the residuals from this regression;
2. Regress Y on X_2, X_3, \dots, X_k , and let \tilde{Y} denote the residuals from this regression;
3. Regress \tilde{Y} on \tilde{X}_1 .

```
# Step 1: regress X1 on X2, X3, ..., Xk, and residuals = X1 tilde
E61c1 <- lm(smoker ~ alcohol + nprevist, data = Birthweight_Smoking)
smoker <- E61c1$residuals

# Step 2: regress Y on X2, X3, ..., Xk, and residuals = Y tilde
E61c2 <- lm(birthweight ~ alcohol + nprevist, data = Birthweight_Smoking)
birthweight <- E61c2$residuals

# Step 3: regress Y tilde on X1 tilde
E61c3 <- lm(birthweight ~ smoker)

summ(E61c3)

## MODEL INFO:
## Observations: 3000
## Dependent Variable: birthweight
## Type: OLS linear regression
##
## MODEL FIT:
## F(1,2998) = 66.55, p = 0.00
## R2 = 0.02
## Adj. R2 = 0.02
##
## Standard errors: OLS
## -----
##               Est.      S.E.    t val.      p
## -----
## (Intercept)      0.00    10.41      0.00     1.00
## smoker          -217.58   26.67     -8.16     0.00
## -----

# comparison
export_summs(E61b, E61c3,
             model.names = c("Model (b)", "Model (c)"))
```

- d. An alternative way to control for prenatal visits is to use the binary variables *Trip0* through *Trip3*. Regress *Birthweight* on *Smoker*, *Alcohol*, *Trip0*, *Trip2*, and *Trip3*.
 - i. Why is *Trip1* excluded from the regression? What would happen if you included it in the regression?
 - ii. The estimated coefficient on *Trip0* is large and negative. What does this coefficient measure? Interpret its value.
 - iii. Interpret the value of the estimated coefficients on *Trip2* and *Trip3*.
 - iv. Does the regression in (d) explain a larger fraction of the variance in birth weight than the regression in (b)?

	Model (b)	Model (c)
(Intercept)	3051.25 *** (34.02)	0.00 (10.41)
smoker	-217.58 *** (26.68)	-217.58 *** (26.67)
alcohol	-30.49 (76.23)	
nprevist	34.07 *** (2.85)	
N	3000	3000
R2	0.07	0.02

*** p < 0.001; ** p < 0.01; * p < 0.05.

Solution

```
E61d <- lm(birthweight ~ smoker + alcohol + tripre0 + tripre2 + tripre3, data = Birthweight_Smoking)
summ(E61d, confint = TRUE, digits = 4)
```

```
## MODEL INFO:
## Observations: 3000
## Dependent Variable: birthweight
## Type: OLS linear regression
##
## MODEL FIT:
## F(5,2994) = 29.1795, p = 0.0000
## R2 = 0.0465
## Adj. R2 = 0.0449
##
## Standard errors: OLS
## -----
##          Est.      2.5%      97.5%      t val.      p
## -----
## (Intercept)    3454.5493    3429.7449    3479.3538    273.0768    0.0000
## smoker         -228.8476    -282.1112    -175.5840     -8.4244    0.0000
## alcohol         -15.1000    -167.1382     136.9383     -0.1947    0.8456
## tripre0        -697.9687    -907.5260    -488.4114     -6.5307    0.0000
## tripre2       -100.8373    -158.9127     -42.7618     -3.4045    0.0007
## tripre3       -136.9553    -253.7798     -20.1308     -2.2986    0.0216
## -----
```

- i. *Tripre1* is omitted to avoid perfect multicollinearity. If we include it in the regression, then coefficient of *Tripre3* disappears.

```
E61di <- lm(birthweight ~ smoker + alcohol + tripre0 + tripre1 + tripre2 + tripre3, data = Birthweight_9)
summ(E61di, confint = TRUE, digits = 4)
```

```
## MODEL INFO:
## Observations: 3000
## Dependent Variable: birthweight
## Type: OLS linear regression
##
## MODEL FIT:
## F(5,2994) = 29.1795, p = 0.0000
## R2 = 0.0465
## Adj. R2 = 0.0449
##
## Standard errors: OLS
## -----
##              Est.          2.5%          97.5%          t val.          p
## -----
## (Intercept)      3317.5941      3201.9115      3433.2766      56.2314      0.0000
## smoker            -228.8476      -282.1112      -175.5840      -8.4244      0.0000
## alcohol           -15.1000      -167.1382       136.9383      -0.1947      0.8456
## tripre0           -561.0135      -798.0225     -324.0044      -4.6412      0.0000
## tripre1             136.9553        20.1308       253.7798       2.2986      0.0216
## tripre2             36.1180       -89.7106       161.9466       0.5628      0.5736
## tripre3
## -----
```

- ii. On average, babies born to women who had no prenatal doctor visits ($Tripre0 = 1$) had birthweights that were 697.9687 grams lower than babies from others who saw a doctor during the first trimester ($Tripre1 = 1$).
- iii.
 - On average, babies born to women whose first doctor visit was during the second trimester ($Tripre2 = 1$) had birthweights that were 100.8373 grams lower than babies from others who saw a doctor during the first trimester ($Tripre1 = 1$).
 - On average, babies born to women whose first doctor visit was during the third trimester ($Tripre3 = 1$) had birthweights that on average were 136.9553 grams lower than babies from others who saw a doctor during the first trimester ($Tripre1 = 1$).
- iv. No, it doesn't. R^2 in (d) is smaller than (b).

```
export_summs(E61b, E61d,
              model.names = c("Model (b)", "Model (d)"))
```

	Model (b)	Model (d)
(Intercept)	3051.25 *** (34.02)	3454.55 *** (12.65)
smoker	-217.58 *** (26.68)	-228.85 *** (27.16)
alcohol	-30.49 (76.23)	-15.10 (77.54)
nprevist	34.07 *** (2.85)	
tripre0		-697.97 *** (106.88)
tripre2		-100.84 *** (29.62)
tripre3		-136.96 * (59.58)
N	3000	3000
R2	0.07	0.05

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$.