

Homework 3: Suggested Solutions

Instructor: Yingyao Hu

By: Tong Zhou

3.22

You estimate a least-squares regression

$$y_i = \mathbf{x}'_{1i} \tilde{\boldsymbol{\beta}}_1 + \tilde{u}_i$$

and then regress the residuals on another set of regressors

$$\tilde{u}_i = \mathbf{x}'_{2i} \tilde{\boldsymbol{\beta}}_2 + \tilde{e}_i.$$

Does this second regression give you the same estimated coefficients as from estimation of a least-squares regression on both set of regressors?

$$y_i = \mathbf{x}'_{1i} \hat{\boldsymbol{\beta}}_1 + \mathbf{x}'_{2i} \hat{\boldsymbol{\beta}}_2 + \hat{e}_i$$

In other words, is it true that $\tilde{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_2$? Explain your reasoning.

Proof. The residual from the regression of \mathbf{y} on \mathbf{X}_1 is:

$$\tilde{\mathbf{u}} = \mathbf{M}_1 \mathbf{y},$$

where $\mathbf{M}_1 = \mathbf{I}_n - \mathbf{P}_1$ and $\mathbf{P}_1 = \mathbf{X}_1(\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1$.

Then $\tilde{\boldsymbol{\beta}}_2$ is obtained by the regression of $\tilde{\mathbf{u}}$ on \mathbf{X}_2 :

$$(1) \quad \tilde{\boldsymbol{\beta}}_2 = (\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2 \tilde{\mathbf{u}} = (\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}.$$

By the *Frisch-Waugh-Lovell Theorem* (or *partialling-out operation*):

$$(2) \quad \hat{\boldsymbol{\beta}}_2 = (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}.$$

Thus, in general $\tilde{\boldsymbol{\beta}}_2$ does not agree with $\hat{\boldsymbol{\beta}}_2$.

There are two extreme conditions under which $\tilde{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_2$:

- \mathbf{X}_1 and \mathbf{X}_2 are orthogonal, i.e. $\mathbf{X}'_1 \mathbf{X}_2 = \mathbf{0}$. Under it, $\tilde{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_2 = (\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{y}$.
- \mathbf{y} is in the column space of \mathbf{X}_1 , i.e. $\mathbf{M}_1 \mathbf{y} = \mathbf{0}$, under which $\tilde{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_2 = \mathbf{0}$. □

3.26

Use the data set from Section 3.22.

- Estimate a log wage regression for the subsample of white male Hispanics. In addition to education, experience, and its square, include a set of binary variables for regions and marital status. For regions, create dummy variables for Northeast, South and West so that Midwest is the excluded group. For marital status, create variables for married, widowed or divorced, and separated, so that single (never married) is the excluded group.
- Repeat this estimation using a different econometric package. Compare your results. Do they agree?

We use *R* and *Stata* to do the regression respectively. [Table 1](#) and [Table 2](#) show that those estimates and standard errors are roughly equal. The only substantial blip is the significance level of the variable `d_NE`. In *Stata*, it is insignificantly different from 0, while in *R* it is significantly different from 0 at the level of 0.1. The reason is that its standard error is around the cut-off point between two significance levels. As a result, rounding errors lead to different inference about its parameter.

Note that some of you did not correctly transform the outcome variable “wage”, which is defined to be $\ln(\text{earnings}/(\text{hours} \times \text{week}))$. See page 94 of Hanse’s book.

Table 1: Regression Results Using Stata

	(1)
	wage
education	0.088*** (0.003)
exp_1	0.028*** (0.003)
exp_2	-0.036*** (0.005)
d_NE	0.062 (0.038)
d_S	-0.068** (0.031)
d_W	0.020 (0.031)
d_married	0.178*** (0.024)
d_WD	0.086** (0.042)
d_sep	0.017 (0.058)
Constant	1.193*** (0.051)
Observations	4230

Standard errors in parentheses

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Table 2: Regression Results Using R

<i>Dependent variable:</i>	
	log_inc
edu	0.087*** (0.003)
exp_1	0.028*** (0.003)
exp_2	−0.036*** (0.005)
d_NE	0.063* (0.038)
d_S	−0.066** (0.031)
d_W	0.018 (0.030)
d_married	0.191*** (0.022)
d_widow	0.091** (0.041)
d_sep	0.020 (0.058)
Constant	1.208*** (0.051)
Observations	4,230

Note: *p<0.1; **p<0.05; ***p<0.01

R codes

```
# load data
data <- read.csv("cps09mar.csv",header = TRUE,sep=",")

# Transform raw variables: log(income), education, experience and
# squared experience
log_inc <- log(data$earnings/(data$hours*data$week))
edu <- data$education
exp_1 <- data$age - data$education - 6
exp_2 <- exp1^2/100

# dummies for regions
d_NE <- ifelse(data$region==1, 1, 0)
d_S <- ifelse(data$region==3, 1, 0)
d_W <- ifelse(data$region==4, 1, 0)

# dummies for marital status
d_married <- ifelse(data$marital==1 | data$marital==2, 1, 0)
d_widow <- ifelse(data$marital==4 | data$marital==5, 1, 0)
d_sep <- ifelse(data$marital==6, 1, 0)

# subsample of white male Hispanics
subsample <- data$race==1 & data$female==0 & data$hispanic == 1

# Run regression
regression <- lm(log_inc ~ edu + exp_1 + exp_2 + d_NE+d_S+d_W+d_
  married+d_widow+d_sep,data=subsample)

# Report
stargazer(regression, title="Regression Results Using R", align=T)
```

Stata codes

```
** 3.26
gen wage = log(earnings/(hours*week))

gen exp_1 = age - education - 6

gen exp_2 = exp_1^2/100

gen d_NE = (region == 1)

gen d_S = (region == 3)

gen d_W = (region == 4 )

gen d_married = (marital == 1 | marital == 2 | marital == 3 )

gen d_WD = ( marital == 4 | marital == 5 )

gen d_sep = ( marital == 6 )

reg wage education exp_1 exp_2 d_NE d_S d_W d_married d_WD d_sep if
    race == 1 & female ==0 & hisp == 1

est store A

esttab A using example5.tex ,b(3) star(* 0.1 ** 0.05 *** 0.01) se(3)
    compress label ///
title(Regression Results Using Stata)
```