

# Solutions for Homework 4

TONG ZHOU

PROBLEMS: 4.16, 4.18, 4.23, 4.25, 4.26, 5.5.

## PROBLEM 4.16

PROOF. 1. Define a composite error

$$v_i = e_i + u_i,$$

then

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + v_i$$

and the assumption  $\mathbb{E}[v_i | \mathbf{x}_i] = 0$  is still maintained.

2. Regress  $y$  on  $\mathbf{x}$ , the estimator is obtained:

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

From  $\mathbb{E}[v_i | \mathbf{x}_i] = 0$ ,

$$\mathbb{E}[\widehat{\boldsymbol{\beta}} | \mathbf{X}] = \boldsymbol{\beta}.$$

So the estimator is unbiased for  $\boldsymbol{\beta}$  despite the presence of measurement error in  $\mathbf{y}$ .

3. The variance of  $\widehat{\boldsymbol{\beta}}$  is:

$$\begin{aligned} \text{Var}(\widehat{\boldsymbol{\beta}} | \mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \underbrace{\mathbb{E}[v v' | \mathbf{X}]}_{=\sigma^2 I_n + \boldsymbol{\Omega}_u} \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \boldsymbol{\Omega}_u \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

where  $\boldsymbol{\Omega}_u = \text{diag}(\sigma_u^2(x_1), \dots, \sigma_u^2(x_n))$ .

Thus, the presence of measurement error in  $\mathbf{y}$  does not affect unbiasedness but inflates variances.  $\square$

## PROBLEM 4.18

PROOF. Define the short regression as

$$\mathbf{y} = \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{e}_1,$$

where  $\mathbf{y} = (y_1 \ y_2 \ \cdots \ y_n)'$ ,  $\mathbf{X}_1 = (X_1 \ \cdots \ X_n)'$ ,  $\mathbf{e}_1 = (e_{11} \ \cdots \ e_{1n})'$ . Here the dimensions  $\mathbf{X}_1$  and  $\mathbf{e}_1$  are  $n \times k_1$  and  $n \times 1$  respectively. The same notations apply to  $\mathbf{X}_2$  and  $\mathbf{y}$ .

Its residual

$$\widehat{\mathbf{e}}_1 = \mathbf{M}_1 \mathbf{e}_1,$$

where  $\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_1(\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1'$  and  $\mathbf{e}_1 = \mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{e}$ .

Thus we have

$$\mathbb{E}[s^2 \mid \mathbf{X}] = \frac{1}{n - k_1} \mathbb{E}[\mathbf{e}_1' \mathbf{M}_1 \mathbf{e}_1 \mid \mathbf{X}].$$

Now it suffices to look into the term  $\mathbb{E}[\mathbf{e}_1' \mathbf{M}_1 \mathbf{e}_1 \mid \mathbf{X}]$ :

$$\begin{aligned} \mathbb{E}[\mathbf{e}_1' \mathbf{M}_1 \mathbf{e}_1 \mid \mathbf{X}] &= \mathbb{E}[(\mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{e})' \mathbf{M}_1 (\mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{e}) \mid \mathbf{X}] \\ &= \mathbb{E}[\boldsymbol{\beta}_2' \mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{e}' \mathbf{M}_1 \mathbf{e} + 2 \boldsymbol{\beta}_2' \mathbf{X}_2' \mathbf{M}_1 \mathbf{e} \mid \mathbf{X}] \\ &= \boldsymbol{\beta}_2' \mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbb{E}[\mathbf{e}' \mathbf{M}_1 \mathbf{e} \mid \mathbf{X}] + 2 \boldsymbol{\beta}_2' \mathbf{X}_2' \mathbf{M}_1 \underbrace{\mathbb{E}[\mathbf{e} \mid \mathbf{X}]}_{=0} \\ &= \boldsymbol{\beta}_2' \mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2 + (n - k_1) \sigma^2. \end{aligned}$$

As a result, we have

$$(1) \quad \mathbb{E}[s^2 \mid \mathbf{X}] = \sigma^2 + \frac{1}{n - k_1} \boldsymbol{\beta}_2' \mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2.$$

Equation (1) implies that  $s^2$  in general is an upward biased estimator for  $\sigma^2$ , since  $\boldsymbol{\beta}_2' \mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2 \boldsymbol{\beta}_2$  is non-negative.  $\square$

#### PROBLEM 4.23

PROOF. Since

$$\mathbb{E}[\widehat{\boldsymbol{\beta}} \mid \mathbf{X}] = (\mathbf{X}' \mathbf{X} + \mathbf{I}_k \lambda)^{-1} \mathbf{X}' \mathbf{X} \boldsymbol{\beta} \neq \boldsymbol{\beta},$$

it is therefore biased for  $\boldsymbol{\beta}$ .  $\square$

#### PROBLEM 4.25

See below.

#### PROBLEM 4.26

See below.

**PROBLEM 5.5**

PROOF. From Hansen's Theorem 5.4,

$$\widehat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}).$$

Then conditional on  $\mathbf{X}$ ,  $\mathbf{x}'_i$  becomes constant for each  $i$ . By theorem 5.2, it follows that  $\widehat{y}_i = \mathbf{x}'_i \widehat{\boldsymbol{\beta}}$  follows a normal distribution such that

$$\widehat{y}_i \mid \mathbf{X} \sim \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}, \sigma^2 h_{ii}).$$

□

## Codes for Problems 4.25 and 4.26

```
rm(list=ls())
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.0 --

## v ggplot2 3.3.3      v purrr 0.3.4
## v tibble 3.1.0       v dplyr 1.0.4
## v tidyr 1.1.2        v stringr 1.4.0
## v readr 1.4.0        v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()

library(knitr)
library(gtsummary)    ## Summarize regression results
library(moderndivde)  ## Summarize regression results
library(haven)        ## Use function read_stata to load dataset.
library(jtools)       ## Report robust standard errors.
library(huxtable)

##
## Attaching package: 'huxtable'

## The following object is masked from 'package:gtsummary':
##
##   as_flextable

## The following object is masked from 'package:dplyr':
##
##   add_rownames

## The following object is masked from 'package:ggplot2':
##
##   theme_grey
```

### Problem 4.25

```
df <- read_stata("~/Dropbox/2021_Spring/Econometrics/hw3/cps09mar.dta")
#### create new variables: log_wage, experience, experience2
df <- df %>% mutate(log_wage = log(earnings/(hours*week)),
                  exper = age - education - 6,
                  exper2 = exper^2/100)

#### create new dummy variables
df <- df %>%
  mutate(d_NE = ifelse(region==1,1,0),
```

```

d_S = ifelse(region==3,1,0),
d_W = ifelse(region==4,1,0),
d_married = ifelse(marital==1 | marital==2,1,0),
d_widow= ifelse(marital==4 | marital==5,1,0),
d_sep = ifelse(marital==6,1,0)) %>%
select(-age)

```

Report the HC3 standard error:

```

#### do regression and obtain table
df %>% filter(race==1,female==0,hisp==1) %>%
  lm(log_wage~education+exper+exper2+d_NE+d_S+
    d_W+d_married+d_widow+d_sep, data=.) %>%
  export_summs(.robust="HC3",digits=3,model.names="Problem 4.25, HC3")

```

## Problem 4.26

Load data and standardize the score variable:

```

ddk <- read_stata("~/Dropbox/2021_Spring/Econometrics/hw4/ddk2011.dta")

ddk <- ddk %>% mutate(std_score = scale(totalscore))

```

Two regressions: conventional s.e. and clustered s.e.

```

model1 <- ddk %>% lm(std_score~tracking+etpteacher+agetest+girl+percentile,data=.)
model2 <- ddk %>% lm(std_score~tracking+etpteacher+agetest+girl+percentile,data=.)

```

Report two regression results

```

export_summs(model1,model2,robust=c("HC3","HC3"),cluster=c("NULL","schoolid"),
  digits=3,model.names = c("Conventional s.e.,"Clustered s.e.))

```

	Problem 4.25, HC3
(Intercept)	1.208 *** (0.051)
education	0.087 *** (0.003)
exper	0.028 *** (0.003)
exper2	-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)
N	4230
R2	0.252

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

	Conventional s.e.	Clustered s.e.
(Intercept)	-0.729 *** (0.081)	-0.729 *** (0.132)
tracking	0.173 *** (0.024)	0.173 * (0.077)
etpteacher	0.180 *** (0.024)	0.180 *** (0.038)
agetest	-0.041 *** (0.009)	-0.041 ** (0.014)
girl	0.081 *** (0.024)	0.081 ** (0.029)
percentile	0.017 *** (0.000)	0.017 *** (0.001)
N	5269	5269
R2	0.249	0.249

Standard errors are heteroskedasticity robust. \*\*\*  $p < 0.001$ ; \*\*  $p < 0.01$ ; \*  $p < 0.05$ .