

## Homework 4: Suggested Solutions

Instructor: Yingyao Hu

By: Tong Zhou

**4.16**

Take the linear homoskedastic CEF

$$(1) \quad \begin{aligned} y_i^* &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}[e_i | \mathbf{x}_i] &= 0 \\ \mathbb{E}[e_i^2 | \mathbf{x}_i] &= \sigma^2 \end{aligned}$$

and suppose that  $y_i^*$  is measured with error. Instead of  $y_i^*$ , we observe  $y_i$  which satisfies

$$y_i = y_i^* + u_i$$

where  $u_i$  is measurement error. Suppose that  $e_i$  and  $u_i$  are independent and

$$\begin{aligned} \mathbb{E}[u_i | \mathbf{x}_i] &= 0; \\ \mathbb{E}[u_i^2 | \mathbf{x}_i] &= \sigma_u^2(\mathbf{x}_i) \end{aligned}$$

1. Derive an equation for  $y_i$  as a function of  $\mathbf{x}_i$ . Be explicit to write the error term as a function of the structural errors  $e_i$  and  $u_i$ . What is the effect of this measurement error on the model [Eq. \(1\)](#)
2. Describe the effect of this measurement error on OLS estimation of  $\boldsymbol{\beta}$  in the feasible regression of the observed  $y_i$  on  $\mathbf{x}_i$ .
3. Describe the effect (if any) of this measurement error on standard error calculation for  $\hat{\boldsymbol{\beta}}$ .

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  $x$ , the estimator is obtained:

$$\hat{\beta} = (X'X)^{-1}X'y$$

From  $\mathbb{E}[v_i|x_i] = 0$ ,

$$\mathbb{E}[\hat{\beta}|X] = \beta.$$

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

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

$$\begin{aligned}\text{Var}(\hat{\beta}|X) &= (X'X)^{-1}X' \underbrace{\mathbb{E}[vv'|X]}_{=\sigma^2 I_n + \Omega_u} X (X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1} + (X'X)^{-1}X'\Omega_u X (X'X)^{-1}\end{aligned}$$

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

Thus, the presence of measurement error in  $y$  does not affect unbiasedness but inflates variances. □

## 4.20

Take the model

$$y = X\beta + e$$

$$\mathbb{E}[e|X] = \mathbf{0}$$

$$\mathbb{E}[ee'|X] = \Omega.$$

Assume for simplicity that  $\Omega$  is known. Consider the OLS and GLS estimators  $\hat{\beta} = (X'X)^{-1}(X'y)$  and  $\tilde{\beta} = (X'\Omega^{-1}X)^{-1}(X'\Omega^{-1}y)$ . Compute the (conditional) covariance between  $\hat{\beta}$  and  $\tilde{\beta}$ :

$$\mathbb{E}[(\hat{\beta} - \beta)(\tilde{\beta} - \beta)'|X]$$

Find the (conditional) covariance matrix for  $\hat{\beta} - \tilde{\beta}$ :

$$\mathbb{E}[(\hat{\beta} - \tilde{\beta})(\hat{\beta} - \tilde{\beta})'|X].$$

Proof. By

$$(2) \quad \hat{\beta} - \beta = (X'X)^{-1}X'e$$

$$(3) \quad \tilde{\beta} - \beta = (X'\Omega^{-1}X)^{-1}X'\Omega^{-1}e,$$

we have

$$\begin{aligned} \mathbb{E}[(\hat{\beta} - \beta)(\tilde{\beta} - \beta)' | X] &= \mathbb{E}[(X'X)^{-1}X'ee'\Omega^{-1}X(X'\Omega^{-1}X)^{-1} | X] \\ &= (X'X)^{-1}X' \underbrace{\mathbb{E}[ee' | X]}_{= \Omega} \Omega^{-1}X(X'\Omega^{-1}X)^{-1} \\ &= (X'X)^{-1}X'X(X'\Omega^{-1}X)^{-1} \\ &= (X'\Omega^{-1}X)^{-1} \end{aligned}$$

For the second part,

$$\begin{aligned} \mathbb{E}[(\hat{\beta} - \tilde{\beta})(\hat{\beta} - \tilde{\beta})' | X] &= \mathbb{E}[(\hat{\beta} - \beta + \beta - \tilde{\beta})(\hat{\beta} - \beta + \beta - \tilde{\beta})' | X]. \\ &= \underbrace{\mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)' | X]}_{= \text{Var}(\hat{\beta} | X)} - \mathbb{E}[(\hat{\beta} - \beta)(\tilde{\beta} - \beta)' | X] \\ &\quad - \mathbb{E}[(\tilde{\beta} - \beta)(\hat{\beta} - \beta)' | X] + \underbrace{\mathbb{E}[(\tilde{\beta} - \beta)(\tilde{\beta} - \beta)' | X]}_{= \text{Var}(\tilde{\beta} | X)} \\ &= (X'X)^{-1}X'\Omega X(X'X)^{-1} - (X'\Omega^{-1}X)^{-1} \\ &\quad - (X'\Omega^{-1}X)^{-1} + (X'\Omega^{-1}X)^{-1} \\ &= (X'X)^{-1}X'\Omega X(X'X)^{-1} - (X'\Omega^{-1}X)^{-1} \end{aligned}$$

□

## 4.23

Take the linear regression model with  $\mathbb{E}[y|X] = X\beta$ . Define the *ridge regression* estimator

$$\hat{\beta} = (X'X + I_k\lambda)^{-1}X'y$$

where  $\lambda > 0$  is a fixed constant. Find  $\mathbb{E}[\hat{\beta} | X]$ . Is  $\hat{\beta}$  biased for  $\beta$ ?

Proof. Since

$$\mathbb{E}[\hat{\beta} | X] = (X'X + I_k\lambda)^{-1}X'X\beta \neq \beta,$$

it is therefore biased for  $\beta$ .

□

**4.26**

Extend the empirical analysis reported in Section 4.22. Do a regression of standardized test score (*totalscore* normalized to have zero mean and variance 1) on tracking, age, gender, being assigned to the contract teacher, and student's percentile in the initial distribution. (The sample size will be smaller as some observation have missing variables.) Calculate standard errors using both the conventional robust formula, and clustering based on the school.

1. Compare the two sets of standard errors. Which standard error changes the most by clustering? Which changes the least?
2. How does the coefficient on *tracking* change by inclusion of the individual controls (in comparison to the results from (4.54))?

See [Table 1](#) for comparison of the clustered s.e. and the conventional HC1 s.e., from which *tracking* changes the most and *gender* changes the least. As a consequence, the p-value of the estimate for *tracking* was 0.026 when the clustered standard errors were used, implying that it was insignificantly different from 0 under 1% significance level.

At the same time, the estimate for *tracking* was 0.174, while it was 0.138 in (4.54). More striking was the change in p-value.

Table 1: Comparison of clustered s.e. and HC1 s.e.

	Coefficients	Clustered s.e.	HC1 s.e.	change by $\Delta\%$
Intercept	-0.742	0.131	0.082	67.8%
tracking	0.174	0.077	0.024	226.5%
assigned	0.181	0.038	0.024	65.5%
age	-0.041	0.013	0.009	23.5%
gender	0.081	0.029	0.024	18.9%
percentile	0.017	0.0007	0.0004	177.6%

(See R codes on the next page)

**R codes:**

```
library(scales)
library(stargazer)

ddk <- read.csv("DDK2011.csv",header= TRUE, sep=",")

## For non-numeric variables, do not forget as.character(),
## otherwise as.numeric() would pump up wrong numbers.

testscore <- as.matrix(ddk$totalscore)
tracking <- as.numeric(as.character(ddk$tracking))
etpteacher <- as.numeric(as.character(ddk$etpteacher))
agetest <- as.numeric(as.character(ddk$agetest))
girl <- as.numeric(as.character(ddk$girl))
percentile <- as.numeric(as.character(ddk$percentile))
schoolid <- as.matrix(ddk$schoolid)

ddk2 <- na.omit(cbind(matrix(1,n,1),tracking,etpteacher,agetest,
  girl,percentile,testscore,schoolid))

y <- scale(ddk2[,7])
n <- nrow(y)

x <- ddk2[,1:6]

schoolid <- ddk2[,8]

k <- ncol(x)

xx <- t(x)%*%x

invx <- solve(xx)

beta <- solve(xx,t(x)%*%y)

## Clustered s.e.

xe <- x*rep(y-x)%*%beta,times=k)
xe_sum <- rowsum(xe,schoolid)
G <- nrow(xe_sum)
omega <- t(xe_sum)%*%xe_sum
scale <- G/(G-1)*(n-1)/(n-k)
V_clustered <- scale*invx)%*%omega)%*%invx
```

```
se_clustered <- sqrt(diag(V_clustered))

## Conventional HC1 s.e.

e = y-x%%beta;
a <- n/(n-k)
u1 <- x*(e%%matrix(1,1,k))

v1a <- a * invx %% (t(u1)%%u1) %% invx
s1a <- sqrt(diag(v1a))

change_by <- percent((se_clustered - s1a)/s1a)

compare <- cbind(beta,se_clustered, s1a,change_by)

colnames(compare) <- c("Coefficients","Clustered s.e.", "HC1 s.e.,"
  change by")

rownames(compare) <- c("Intercept","tracking","assigned","age","
  gender","percentile")

stargazer(compare)
```