# Econ 613 - Applied Econometrics - 2022 Spring
# Homework 4

### Yican Liu

### April 20, 2022

Before I do the analysis, I dropped the observations with NA or negative income.

## 1   Exercise 1

### 1.1

I use the current year (2019) minus the birth year (KEY _BDATE_Y_1997) to measure the age. I aggregate all "CV_WKSWK_JOB_DLI" to create the total work experience.
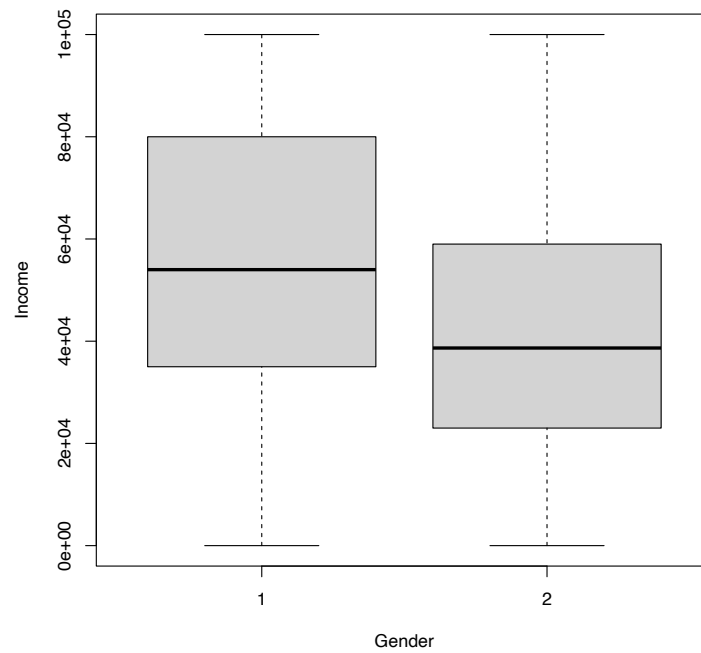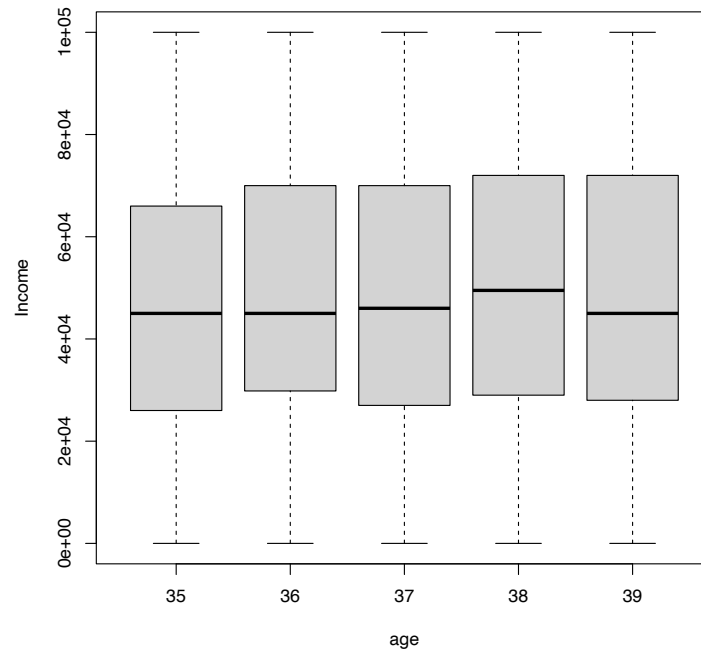
### 1.2

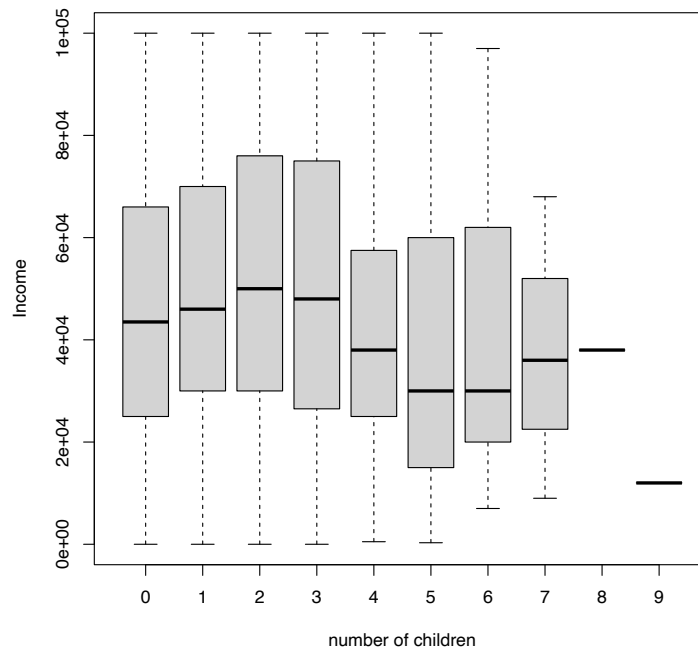I transfer YSCH.3113_2019 to numerics. The schooling year is transfered as follows:

- GED: 4 years.

- High school: 12 years.

- Bachelor: 16 years.

- Master: 18 years.

- Ph.D.: 23 years.

- Professional: 22 years.

- AA: 14 years.

Other variables, such as those for parents' education, are not transferred because these variables indicate the total education year. I generate a new variable, *total_edu*, as the transferred variable of YSCH.3113_2019. I use *total_edu* in the following analysis.

## 1.3

### 1.3.1 Income by Groups

### 1.3.2 The Share of Zero in Income

The proportion of zero income by age, gender, marital status, and number of children is shown as follows.

```
  age % of zero income
1  35       0.009293680
2  36       0.006300630
3  37       0.005420054
4  38       0.008960573
5  39       0.002994012
```

```
  gender % of zero income
1      0       0.007500000
2      1       0.005742726
```

```
  marital status % of zero income
1              0       0.005592272
2              1       0.007454342
3              2       0.043010753
4              3       0.001538462
5              4       0.000000000
```

```
       number of children % of zero income
1                        0        0.006993007
2                        1        0.007846556
3                        2        0.005743001
4                        3        0.008025682
5                        4        0.000000000
6                        5        0.000000000
7                        6        0.000000000
8                        7        0.000000000
9                        8        0.000000000
10                       9        0.000000000
```

### 1.3.3

To interpret these results, we know the following facts.

- In general, the average incomes of different age groups (from 35 to 39) are almost the same. The age group of 36 has a slightly higher average income compared to the other age groups. The zero income rate also does not vary a lot across all these five groups.

- Male has a higher average income than female. In addition, the female has a higher zero income rate compared to male.

- People who have one, two, three, or four children tend to earn more than the other people. In addition, those who have no children have the highest zero income rate.

## 2  Exercise 2

### 2.1

### 2.1.1

The OLS regression is shown as follows

$$\text{Income}_i = \beta_0 + \beta_1 \text{Age}_i + \beta_2 \text{WorkExp}_i + \beta_3 \text{Education}_i + \beta_4 \text{Female}_i + \beta_5 \text{Child}_i + \epsilon_i \qquad (1)$$

where $i$ indicates the observation (individual). I include age, work experience, education, number of children, as well as female dummy variable in my model. The regression coefficients are as follows:

```
Call:
lm(formula = data.data4.posinc$YINC_1700_2019 ~ data.data4.posinc$age +
    data.data4.posinc$work_exp + data.data4.posinc$total_edu +
    data.data4.posinc$children + data.data4.posinc$female)

Residuals:
   Min     1Q Median    3Q    Max
-70981 -18433  -3724  16354  89175

Coefficients:
                             Estimate Std. Error t value Pr(>|t|)
(Intercept)                  12765.21    9462.54   1.349    0.177
data.data4.posinc$age          331.17     254.52   1.301    0.193
data.data4.posinc$work_exp     999.16      66.26  15.080  < 2e-16 ***
data.data4.posinc$total_edu   1554.28      76.14  20.413  < 2e-16 ***
data.data4.posinc$children    1323.13     287.47   4.603 4.27e-06 ***
data.data4.posinc$female    -15185.78     722.33 -21.023  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 25920 on 5370 degrees of freedom
Multiple R-squared:  0.1728,    Adjusted R-squared:  0.172
F-statistic: 224.4 on 5 and 5370 DF,  p-value: < 2.2e-16
```

In this regression, the regression coefficients of work experience, number of children, work experience, and female are significant. However, the regression coefficient for age is not significant. This is not commonsensible. Thus, there might be biases when estimating the regression coefficient by OLS regression.

To interpret the results, for the people who have positive income, we know that: given everything else fixed, if work experience increases by 1 unit, the personal income will also increase by 999.16 dollars; if personal total education increases by 1 year, the personal income will also increase by 1554.28; if the number of children increases by 1, the personal income will also increase by 1323.13 dollars; and for female individuals, given anything else as controlled, they earn 15185.78 less than male individuals.

### 2.1.2

Selection problem exists in this model because we only focus on the people who have positive income. That is, we dropped the observations with zero income.

### 2.2

In our cases, the endogeneity problem exists due to selection biases. Thus, by using Heckman methodology, we include a new variable $\text{IMR}_i$ in our regression analysis to control for the probability of "having a positive income."

### 2.3

I estimate the Heckman Model as follows.

First, I run a Probit model to estimate the probability for whether $\text{Income}_i > 0$.

Second, for each observation $i$, I estimate the value of

$$\text{IMR}_i = \frac{\phi\left(\frac{x_i\beta}{\sigma}\right)}{\Phi\left(\frac{x_i\beta}{\sigma}\right)}$$

and include this value in the OLS regression (1). The regression coefficient is

```
Call:
lm(formula = data.data4.heckman$YINC_1700_2019 ~ data.data4.heckman$age +
    data.data4.heckman$work_exp + data.data4.heckman$total_edu +
    data.data4.heckman$children + data.data4.heckman$female +
    data.data4.heckman$imr)

Residuals:
    Min      1Q  Median      3Q     Max
 -65625  -18452   -3770   16574   93037

Coefficients:
                                Estimate Std. Error t value Pr(>|t|)
(Intercept)                      141016.8    34382.1   4.101 4.17e-05 ***
data.data4.heckman$age            -2097.1      675.6  -3.104 0.001917 **
data.data4.heckman$work_exp         228.2      209.5   1.089 0.276093
data.data4.heckman$total_edu       1015.0      158.5   6.405 1.63e-10 ***
data.data4.heckman$children        -249.2      496.7  -0.502 0.615901
data.data4.heckman$female        -18893.6     1197.4 -15.779  < 2e-16 ***
data.data4.heckman$imr         -1126484.3   290360.0  -3.880 0.000106 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 25890 on 5369 degrees of freedom
Multiple R-squared:  0.1751,    Adjusted R-squared:  0.1742
F-statistic:   190 on 6 and 5369 DF,  p-value: < 2.2e-16
```
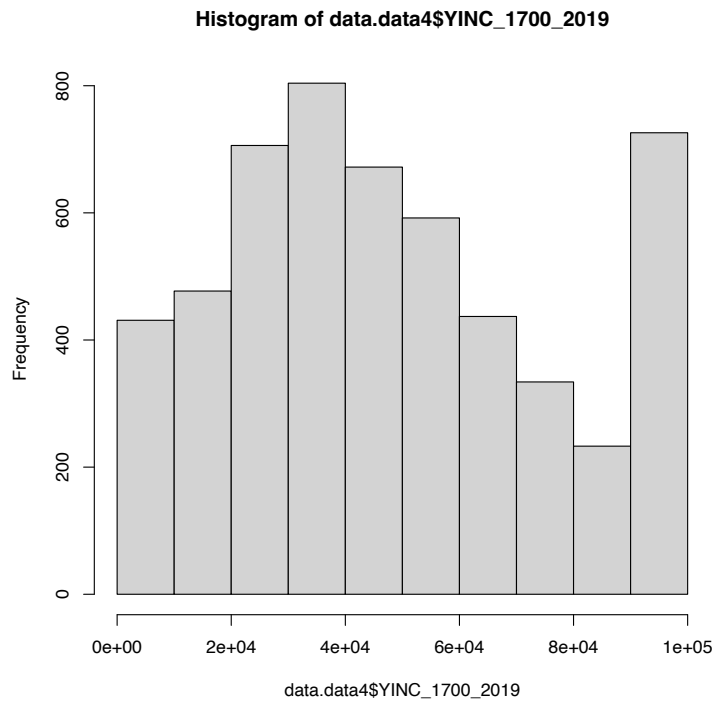
In this table, we can see that in the Heckman model, the regression coefficient for age, personal total education (generated by YSCH.3113_2019), female dummy, as well as a Heckman IMR has significant regression results. Comparing to our previous regression results, we can know that the Heckman IMR variable is an omitted variable in our previous empirical analysis.

In the OLS setting, the regression coefficients for work experience and the number of children are significant. However, in the Heckman regression, these two variables are not significant. Age becomes significant in the Heckman regression, while it is not a significant regressor in the previous OLS regression. Given anything else fixed, if the age increase by 1, the personal income will decrease by 2097.1; if personal total education year increased by 1 year, the income will increase by 1015.0; and given everything else fixed, female individuals earn 18893.6 less than male individuals. After controlling for people's willingness to work and earn money (Heckman IMR), the effects of work experience become less significant.

# 3 Exercise 3

## 3.1

The histogram for income variable is as follows

**Histogram of data.data4$YINC_1700_2019**



As shown in this figure, the censored value might be 100000.00.

### 3.2

I use Tobit model to deal with such problems (the Tobit model for right-hand side censoring).

### 3.3

I use *optim* command in R to estimate the Tobit model regression coefficient. The results are shown as follows:

```
$par
                                     [,1]
(Intercept)                     12761.8020
data.data4.posinc$age             309.7582
data.data4.posinc$work_exp       1078.9516
data.data4.posinc$total_edu      1666.8389
data.data4.posinc$female        -15341.8900
data.data4.posinc$children       1396.4698
log_sigma                          10.2724

$value
[1] 56180.66

$counts
function gradient
     142      100

$convergence
[1] 1
```

## 3.4

To interpret the results above, we know that given everything else fixed, if age increased by 1, the personal income will increase by 309. On average, female individuals earn 15341.89 less than male individuals.

For comparison, I also use the OLS regression to estimate the regression coefficients. The results are shown as follows:

```
Call:
lm(formula = data.data4.posinc$YINC_1700_2019 ~ data.data4.posinc$age +
    data.data4.posinc$work_exp + data.data4.posinc$total_edu +
    data.data4.posinc$female + data.data4.posinc$children)

Residuals:
   Min     1Q Median    3Q    Max
-68563 -20516  -3870  18420  67533

Coefficients:
                              Estimate Std. Error t value Pr(>|t|)
(Intercept)                   36127.08    9763.34   3.700 0.000218 ***
data.data4.posinc$age           274.00     264.20   1.037 0.299743
data.data4.posinc$work_exp     1163.54      68.31  17.034  < 2e-16 ***
data.data4.posinc$total_edu      33.26      24.38   1.364 0.172479
data.data4.posinc$female     -13649.04     745.78 -18.302  < 2e-16 ***
data.data4.posinc$children     1175.44     298.61   3.936 8.38e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 26900 on 5370 degrees of freedom
Multiple R-squared:  0.1089,    Adjusted R-squared:  0.1081
F-statistic: 131.3 on 5 and 5370 DF,  p-value: < 2.2e-16
```

The regression coefficient is slightly different from the Tobit model. This is because we did not consider the variance of income from income above 100000.00.

# 4 Exersice 4

## 4.1

The potential biases come from something that varies from person to person (for example, personality) and something that varies from time to time (for example, macroeconomic conditions), which is unobservable but needs to control. That is, including fixed effects might solve the potential biases.

## 4.2

I reorganize the wide panel data to long panel data before I do these regressions.

### 4.2.1 Within Estimator

I estimate the within estimator as

$$(\text{Income}_{it} - \overline{\text{Income}}_i) = \beta_0 + \beta_1(\text{Edu}_{it} - \overline{\text{Edu}}_i) + \beta_2(\text{WorkExp}_{it} - \overline{\text{WorkExp}}_i) + (\epsilon_{it} - \bar{\epsilon}_i)$$

and the regression coefficient is

```
Call:
lm(formula = data.panel$income.demeaned ~ data.panel$edu.demeaned +
    data.panel$work.demeaned)

Residuals:
    Min      1Q  Median      3Q     Max
-110449   -7356       0    5690  304071

Coefficients:
                          Estimate Std. Error t value Pr(>|t|)
(Intercept)             -7.207e-10  4.401e+01     0.0        1
data.panel$edu.demeaned  1.091e+03  1.051e+01   103.8   <2e-16 ***
data.panel$work.demeaned 6.131e+01  3.552e-01   172.6   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 17700 on 161709 degrees of freedom
Multiple R-squared:  0.2929,    Adjusted R-squared:  0.2929
F-statistic: 3.349e+04 on 2 and 161709 DF,  p-value: < 2.2e-16
```

### 4.2.2   Between Estimator

I estimate the between estimator as

$$\overline{\text{Income}}_i = \beta_0 + \beta_1 \overline{\text{Edu}}_i + \beta_2 \overline{\text{WorkExp}}_i + \bar{\epsilon}_i$$

and the regression coefficient is

```
Call:
lm(formula = data.panel.between$mean.income ~ data.panel.between$mean.total.edu +
    data.panel.between$mean.total.work)

Residuals:
   Min     1Q Median     3Q    Max
-38010  -5725  -1168   3195  87129

Coefficients:
                                      Estimate Std. Error t value Pr(>|t|)
(Intercept)                             12.089    216.500   0.056    0.955
data.panel.between$mean.total.edu     1027.447     26.264  39.120   <2e-16 ***
data.panel.between$mean.total.work      50.620      1.293  39.151   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 10250 on 8981 degrees of freedom
Multiple R-squared:  0.3772,    Adjusted R-squared:  0.3771
F-statistic:  2720 on 2 and 8981 DF,  p-value: < 2.2e-16
```

### 4.2.3   Difference Estimator

I estiamte the difference estimator as

$$(\text{Income}_{it} - \text{Income}_{it-1}) = \beta_0 + \beta_1(\text{Edu}_{it} - \text{Edu}_{it-1}) + \beta_2(\text{WorkExp}_{it} - \text{WorkExp}_{it-1}) + (\epsilon_{it} - \epsilon_{it-1})$$

and the regression coefficient is

```
Call:
lm(formula = data.panel$income.diff ~ data.panel$total.edu.diff +
    data.panel$total.work.diff)

Residuals:
    Min      1Q  Median      3Q     Max
-236253   -3574   -1143    2641  345176

Coefficients:
                            Estimate Std. Error t value Pr(>|t|)
(Intercept)               1142.8981    43.0296   26.56   <2e-16 ***
data.panel$total.edu.diff  761.5424    11.5574   65.89   <2e-16 ***
data.panel$total.work.diff  36.8401     0.4115   89.53   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 16650 on 152725 degrees of freedom
  (8984 observations deleted due to missingness)
Multiple R-squared:  0.1153,    Adjusted R-squared:  0.1153
F-statistic:  9950 on 2 and 152725 DF,  p-value: < 2.2e-16
```