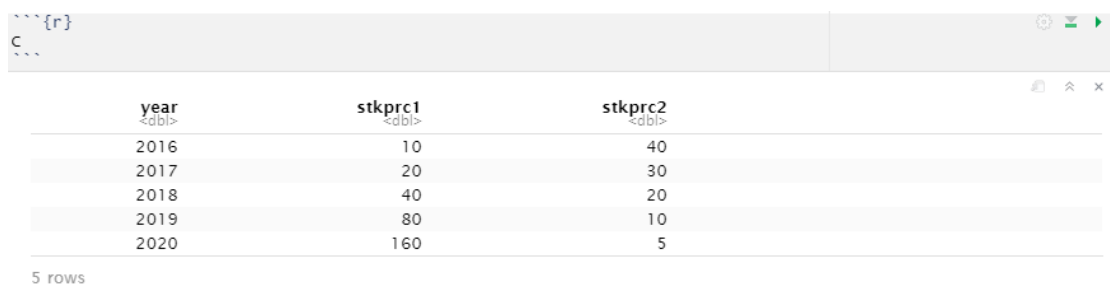# Final Exam of Financial Econometrics, Fall 2022

December 26, 2022
8:30-10:30 am
Total Points: 120

1. (8 points, 2 points each)

    For each of the following R codes, describe what the code is trying to do, and write down the output:

```{r}
C
```

| year <dbl> | stkprc1 <dbl> | stkprc2 <dbl> |
|---|---|---|
| 2016 | 10 | 40 |
| 2017 | 20 | 30 |
| 2018 | 40 | 20 |
| 2019 | 80 | 10 |
| 2020 | 160 | 5 |

5 rows

  a) C[1,]

```
               !
year stkprc1 stkprc2
2016 10      40
```

  b) C[c(1,2),c(1,3)]

```
year stkprc2
2016 40
2017 30
```

  c) C[year>=2019]

```
                2019
year stkprc1 stkprc2
2019 80      10
2020 160     5
```

  d) library(dplyr)

    (C %>% select(year, stkprc2) %>% filter(year<2018))

```
          2018   year stkprc2
year stkprc2
2016 40
2017 30
```

2. (9 points, 3 points each) Write R codes to perform the tasks required.

```r
D
```

| | year <dbl> | stkprc1 <dbl> |
|---|---|---|
| 1 | 2016 | 10 |
| 2 | 2017 | 20 |
| 3 | 2018 | 40 |
| 4 | 2019 | 80 |
| 5 | 2020 | 160 |

5 rows

```r
E
```

| | year <dbl> | stkprc2 <dbl> |
|---|---|---|
| 1 | 2016 | 40 |
| 2 | 2017 | 30 |
| 3 | 2018 | 20 |
| 4 | 2019 | 10 |

4 rows

a) Using D and E, assemble a new data frame F.

F should have three variables: year, stkprc1 and stkprc2.

F should have 5 rows, corresponding to year 2016, 2017, 2018, 2019, 2020.

Your code should take only one line.

Just write down the code. No need to write down the output.

```r
F <- merge(D, E, by = "year", all.x = TRUE)
```

b) Using D and the ggplot2 package, draw a scatter plot of stkprc1 (y) on year (x) with a linear regression line and the associated confidence intervals of the linear regression.

Just write down the code. No need to draw the figure yourself.

```r
ggplot(D, aes(x = year, y = stkprc1)) + geom_point() + geom_smooth(method = "lm", se = TRUE)
```

c) Using C and the ggplot2 package, draw a line plot of stkprc1 (y1) and stkprc2 (y2) on year (x).

Just write down the code. No need to draw the figure yourself.

```r
ggplot(C, aes(x = year)) +
  geom_line(aes(y = stkprc1), color = "blue") +
  geom_line(aes(y = stkprc2), color = "red")
```

3. (6 points) You are studying the relationship between a firm's return on equity (ROE) and CEO salary. In the dataset ceosal, you find that the average salary of a CEO (salary) is 1,281 thousand USD, and the average ROE of a firm (roe) is 17.18%. The covariance between salary and roe is 1342.54. The standard deviation of salary is 1,372 thousand USD, and the standard deviation of roe is 8.52%. The regression equation is:

$$salary\_i = b0 + b1 * roe\_i + u\_i$$

What is the OLS estimate of b1?

```
b1 = cov(salary_i , roe_i)/var(roe_i) = 1342.54 / 0.0852 ^ 2 = 184947.32
[       8.52    8.52%?]
```

4. (5 points) In the linear regression model, we have defined the total sum of squares (TSS), the explained sum of squares (SSE), and the residual sum of squares (SSR). They are calculated as follows:

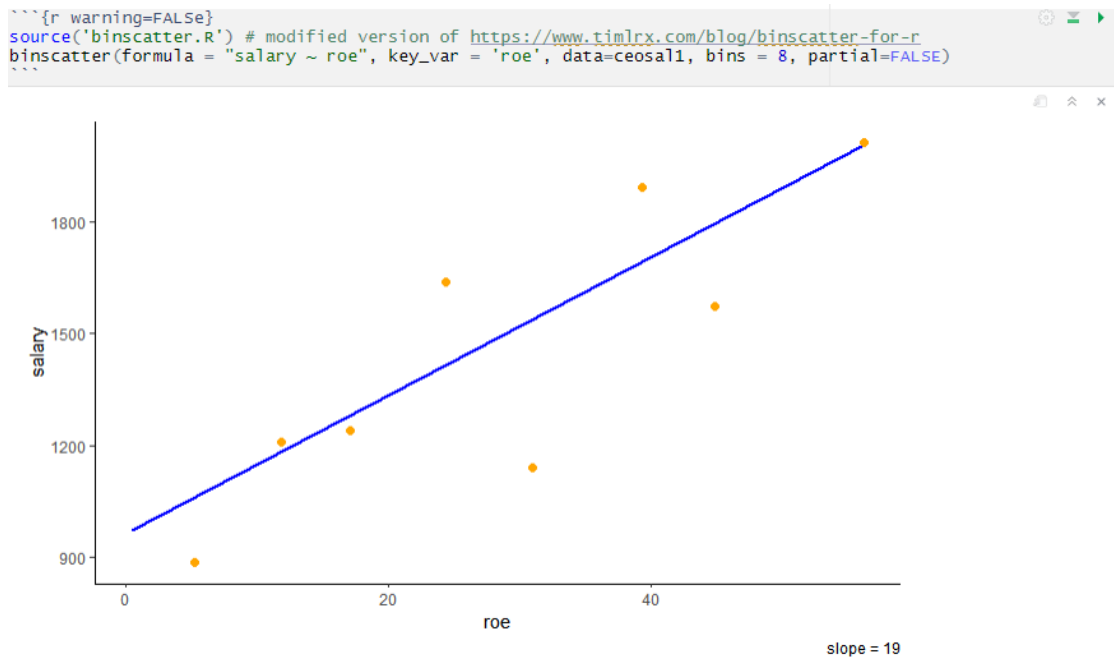$$TSS = \sum_{i=1}^{n} (y_i - \bar{y})^2, SSE = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2, SSR = \sum_{i=1}^{n} \hat{u}_i^2$$

where $\hat{y}_i = b_0 + b_1 x_{1i} + \cdots + + b_k x_{ki}$ and $\hat{u}_i = y_i - \hat{y}_i$.

Question: Which of the following is the R-squared?

(A) $\left(\frac{SSE}{SSR}\right)$      (B) $\left(\frac{SSE}{SSR}\right)^2$      (C) $\left(\frac{SSE}{TSS}\right)$      (D) $1 - \left(\frac{SSE}{TSS}\right)$

R        /

1-SSR/TSS

5. (5 points, 1 point each) True or False (binscatter)

```{r warning=FALSE}
source('binscatter.R') # modified version of https://www.timlrx.com/blog/binscatter-for-r
binscatter(formula = "salary ~ roe", key_var = 'roe', data=ceosal1, bins = 8, partial=FALSE)
```



slope = 19

A) The code generates a scatter plot of salary (y) on roe (x) using the ceosal1 dataset.

B) The code generates a scatter plot of grouped averages of salary (y) on roe (x) using the ceosal1 dataset.

C) The figure indicates that the homotheticity assumption likely does not hold.

D) The figure indicates that the relationship between salary and roe may not be as simple as a linear relationship.

E) The binscatter graph is a helpful tool, but it is not useful for exploring the relationship between y and x when y contains measurement errors.
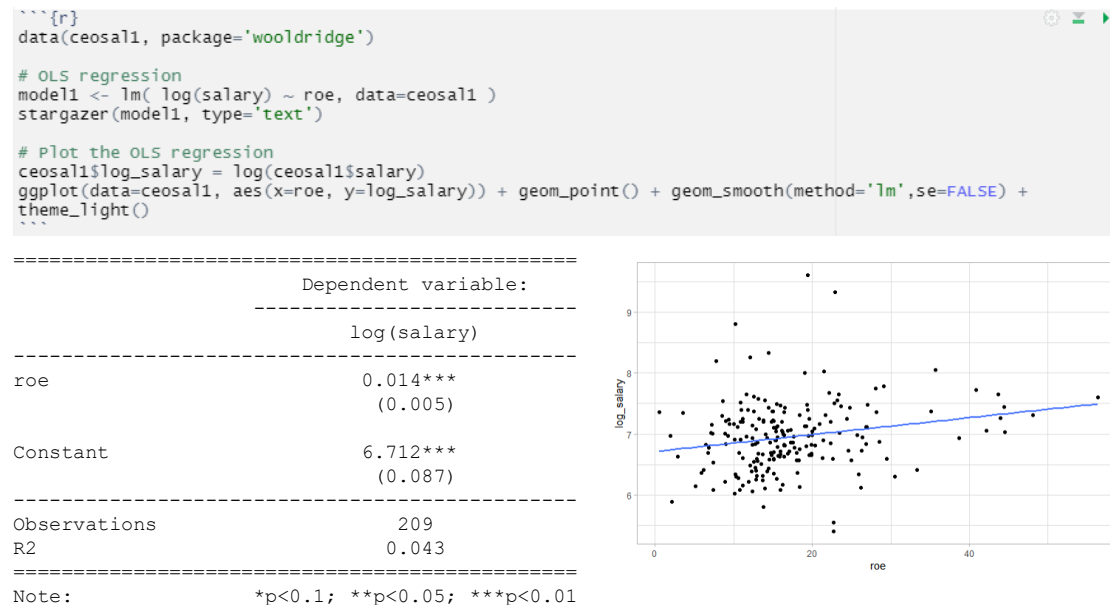
a) False

b) True
bin=8

c) True

d) False

e) False

6. (9 points, 3 point each) Interpret the regression results in the following table. The dependent variable is the logarithm of the salary of a firm's CEO, and the independent variable is the firm's return on equity (ROE).

```r
```{r}
data(ceosal1, package='wooldridge')

# OLS regression
model1 <- lm( log(salary) ~ roe, data=ceosal1 )
stargazer(model1, type='text')

# Plot the OLS regression
ceosal1$log_salary = log(ceosal1$salary)
ggplot(data=ceosal1, aes(x=roe, y=log_salary)) + geom_point() + geom_smooth(method='lm',se=FALSE) +
theme_light()
```
```

```
===============================================
                   Dependent variable:
                  -----------------------------
                        log(salary)
-----------------------------------------------
roe                      0.014***
                         (0.005)

Constant                 6.712***
                         (0.087)
-----------------------------------------------
Observations               209
R2                        0.043
===============================================
Note:                *p<0.1; **p<0.05; ***p<0.01
```



a) First, look at the regression results on the left. Interpret the economic size and the statistical significance of the regression coefficient of log(salary) on roe.

ROE

ROE   log  salary                          t        =0.014/0.005=2.8           p          0.01***

b) Now, look at the figure on the right. Is the standard error in part a) correct? Why?

x                       x

1128      01  17  18
          c)   A friend recommended you to run the following code in addition to the code above.
               What does the code do and what is the meaning of "HC" in "HC0"?

               coeftest(model1, vcov=vcovHC, type=HC0)

coeftest

        vcov = vcovHC            vcovHC

" HC"        " Heteroskedasticity-Consistent"
" HC0"

7. (6 points, 3 point each) Interpreting the following code and result. The data set contains the wage level (wage), the number of years of education (educ), and the number of years of work experience (exper) for 526 workers.

```r
library(car)
data(wage1, package='wooldridge')
model3 <- lm(log(wage) ~ educ+exper, data=wage1)
stargazer(model3, type='text')
linearHypothesis(model3, "educ=exper")
```

```
===============================================
                      Dependent variable:
                    ---------------------------
                            log(wage)
-----------------------------------------------
educ                         0.098***
                              (0.008)

exper                        0.010***
                              (0.002)

Constant                     0.217**
                              (0.109)

-----------------------------------------------
Observations                    526
R2                             0.249
Adjusted R2                    0.246
Residual Std. Error     0.461 (df = 523)
F Statistic          86.862*** (df = 2; 523)
===============================================
Note:                *p<0.1; **p<0.05; ***p<0.01
Linear hypothesis test

Hypothesis:
educ - exper = 0

Model 1: restricted model
Model 2: log(wage) ~ educ + exper

  Res.Df    RSS Df Sum of Sq       F    Pr(>F)
1    524 141.92
2    523 111.34  1    30.576 143.62 < 2.2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

a) Look at the last line of code that starts with "linearHypothesis(…" What does this line of code do?

linearHypothesis()

              educ      exper

b) What can you say about the effects of one more year of education versus one more year of work experience, respectively, on the worker's wage?

education:          ,            0.098

exper           ,            0.010

8. (6 points, 3 points each)

Interpret the effect of x variables on wage:

```
===============================================
                      Dependent variable:
                    ---------------------------
                            log(wage)
-----------------------------------------------
married                      0.213***
                              (0.055)

female                       -0.110**
                              (0.056)

educ                         0.079***
                              (0.007)

exper                        0.027***
                              (0.005)

I(exper2)                    -0.001***
                              (0.0001)

married:female               -0.301***
                              (0.072)

Constant                     0.321***
                              (0.100)

-----------------------------------------------
Observations                    526
R2                             0.461
Adjusted R2                    0.453
Residual Std. Error     0.393 (df = 517)
F Statistic           55.246*** (df = 8; 517)
===============================================
Note:               *p<0.1; **p<0.05; ***p<0.01
```

a) Everything else equal, how does the wage of a single male compare to the wage of a married female?

```
    +      VS      +

 married  female  married  female
         married=0 female=0  : y = constant
         married=1 female=1   y = constant + 0.213 - 0.110 - 0.301 = constant - 0.198
                     0.198
```

b) For a worker with 10 years of work experience, what is the expected effect of an additional year of work experience on her wage?

(Hint: the answer is slightly different if you use an exact calculation or if you take the derivative. Both answers will be correct.)

```
                    exper  I  exper2

  1
exper=10, y1 = constant + 10*0.027 + 100*  -0.001  = constant+0.17
exper=11, y2 = constant + 11*0.027 + 121*(-0.001) = constant +0.176
y2 - y1 = 0.006
     10                                                         0.006

  2
dy/dx = 0.027 - 0.002 * exper
     exper=10     dy/dx =0.007
        10                                                      0.007
```

9.  (10 points, 2 points each) True or False (Binary choice models)

a) After estimating the linear probability model, you should compute the marginal effect.

b) After estimating the logit model or the probit model, you should compute the marginal effect.

c) The marginal effect of an x variable measures the average increase in Prob(Y=1) when that x variable increases by one.

d) After estimating a logit or a probit model, you should check the predicted value of the regression to see whether the predicted values are mostly between 0 and 1.

e) The logit model and the probit model are better than the linear probability model because the effect on Prob(Y=1) when the value of a dummy variable increases from 0 to 1 is more stable under the logit model and the probit model.

Binary choice models            week4

a) FALSE

b) TRUE

marginal effect

logit                                              marginal effect

c) TRUE                 x                  1

d) FALSE LPM      0  1          logit  probit

e) TRUE

10. (6 points) Measurement errors in the x variable: Suppose the true model that generates y and x is as follows:

$$y_i = \alpha + \beta x_i + e_i$$

However, you do not observe the true value $x_i$. Instead, you observe $x_i^* = x_i + v_i$, where $v_i$ is the measurement error in the x variable.　　　　X

Thus, when you estimate the regression model, you are in fact estimating the following model with measurement error in the x variable:

$$y_i = b_0 + b_1 x_i^* + u_i$$

Question: What is the relationship between $\hat{b}_1$ (the estimated regression coefficient of y on x* when there are measurement errors in the x variable) and $\beta$ (the true regression coefficient)?

Describe the steps you took to get to the answer.

```
b1^hat                          b1^hat       0

                                yi  = b0+b1xi +ui +b1vi (     xi * = xi +vi )

            =cov(x  y)/var(x

    b1=cov(yi , xi )/var(xi )=[cov(b0, xi )+cov(b1xi , xi )+cov(ui , xi )+cov(b1v1, xi *-v1)/var(xi )

        =[cov(b0, xi )+cov(b1xi , xi )+cov(ui , xi )]/var(xi )                cov(b1xi , xi )   0

    b1=   + cov(b1v1   xi *-v1)/var(xi )  ==>       cov(b1v1   xi *-v1)/var(xi )


(                                                                    xi = xi *+vi



b1^hat =    · var(xi ) /[ var(xi )+var(vi )]--
```

11. (8 points, 2 points each) True or False (Tobit and Heckman selection)

a) The OLS is not valid when the data is MCAR.

b) The Tobit model is designed to handle situations where you observe a censored version of the true dependent variable. For example, the true dependent variable is y and you observe only y* = max(0,y). If you would like to know the effect of x on y, you should compute the marginal effect after estimating the Tobit model.

c) The Heckman selection model is designed to handle situations where selection into the sample may depend on some unobserved benefit, so that an estimation based on the OLS will suffer from an omitted variable bias. The Heckman selection model includes the Inverse Mills ratio from the second stage estimation as an additional control variable in the first stage regression to address this omitted variable bias concern.

d) You can use exactly the same set of x variables in the first stage and the second stage of the Heckman selection model.

tobit heckman

a) FALSE MCAR=missing completely at random week5 5.1
　　　　 "　　　　　　　　 "　　 MCAR

b) TRUE censored version　 tobit model week5 5.3
　　　　　　　　　　　　　　 censored

c) FALSE heckman　　　　　　 week5 5.4
　　　　　　　　　　　 MLE　　 Probit　　　　　　　　　 inverse mills ratio


　　　　　　　　　　　　　　　　 '

ZY

d) FALSE

12. (9 points, 3 points each) The instrumental variable approach can be described as follows. You are studying the relationship between y and x:

$$y_i = b_0 + b_1 x_i + u_i$$

but $x_i$ and $u_i$ are not independent, which violates the OLS assumptions and biases the estimated coefficient $\hat{b}_1$.

For example, $y_i$ is a firm's innovation activities (e.g. number of new products developed), and $x_i$ is the amount of VC investment the firm receives; in this case, $u_i$ can include many factors that affects the firm's $y_i$ beyond VC investment that it received.

a) Suppose whether the firm's founder graduated from a top 2 university positively attracts VC investment and positively affects innovation activities. How will this affect the correlation between $x_i$ and $u_i$, and how will this affect the estimated coefficient $\hat{b}_1$?

xi ui          b1^hat                    cov(yi , xi )                    cov   ui , xi

b) The instrumental variable approach uses an instrumental variable $z_i$, and involves estimating a first stage regression, computing the predicted value $\hat{x}_i$, and use the predicted value in the original regression:

$$x_i = a_0 + a_1 z_i + u_i$$
$$y_i = b_0 + b_1 \hat{x}_i + u_i$$

What are the names and the contents of the two assumptions required in order for the instrumental variable $z_i$ to be valid?

z       x       cov   xi    zi              0
z       ui            z       x       y

c) Based on these two assumptions, evaluate whether the number of VC funds in the firm's headquarter city is a valid instrument for the above example. Explain why.
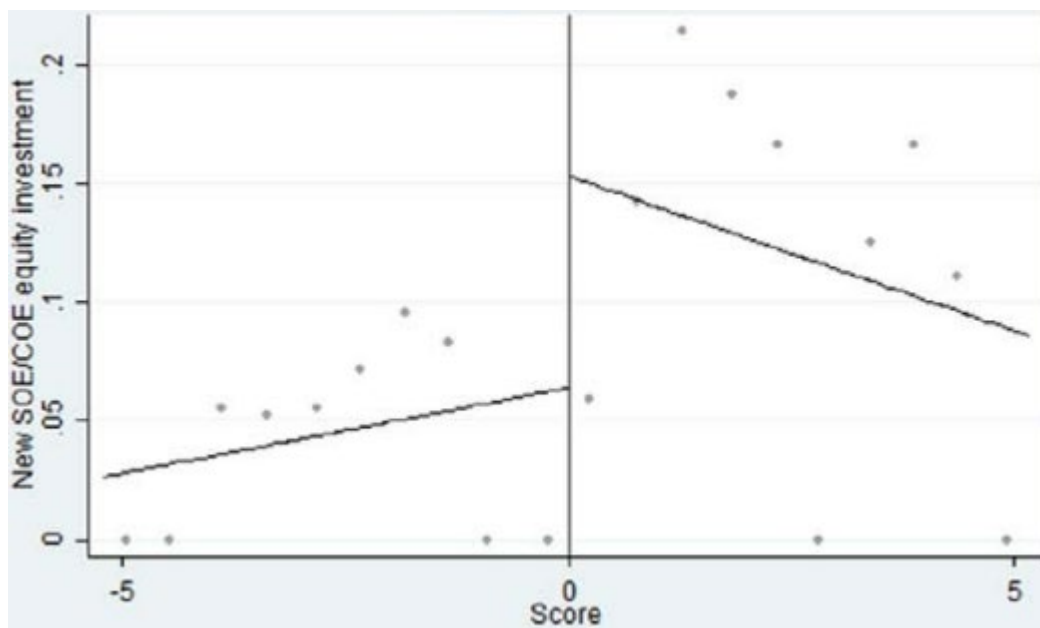
VC

13. (7 points) Write down the regression equation for the following regression discontinuity plot.

The name of the y variable is SOEInvestment (vertical axis), i.e. equity investment the firm received from state-owned investors.

The horizontal axis is a firm's score in the Innofund (科技部创新基金) grant evaluation process.

When Score>=0, the firm successfully gets the grant and certification from the Innofund.

Specify which regression coefficient is the coefficient of interest in the regression equation you write down.

14. (14 points, 2 points each) True or False (Time series and panel data)
   a) You can control for seasonality in a time series regression with monthly data by including month of the year dummies (Jan, Feb, …) as control variables.
   b) The unit root test uses a regression based on the AR(1) model:
$$y_t = \alpha + \rho y_{t-1} + e_t$$
   If you cannot reject the null hypothesis $\rho = 0$, then $y_t$ has a unit root.
   c) When $y_t$ has a unit root, $\Delta y_t = y_t - y_{t-1}$ will NOT have a unit root.
   d) HAC standard errors can deal with regression errors that are correlated over time.
   e) Clustered standard errors can deal with regression errors that are correlated over time, but can only be used with panel data.
   f) With panel data, you can control for unobservable characteristics that are time-invariant using individual fixed effects. Then, the regression coefficient measures the effect of the x variable, holding fixed the same individual.
   g) With panel data, you can control for the effect of macroeconomic and aggregate factors using time fixed effects.

15. (6 points, 2 points each) True or False (Event studies)
   a)   The for-loop version of the event study R code we studied in class explicitly addresses the issue that events are overlapping (so that abnormal returns may be correlated across events) by using bootstrap standard errors, so that the t-stats produced by that R code is correct.
   b)   The for-loop version of the event study R code we studied in class did not address the issue that events are overlapping (so that abnormal returns may be correlated across events.) Therefore, the average CAR produced by that R code is incorrect.
   c)   After defining the event and the window parameters, the event study estimates the normal return and the abnormal returns using data from the event window.

16. (6 points) You are given with the following regression equation:

$$R_{it} = \alpha_i + \beta_i R_{mt} + \sum_{s=T_1}^{s=T_2} \gamma_s 1\{t = s\} + \epsilon_{it}$$

   a)   Please explain the meaning and role of each of the variables and the coefficients in this regression equation. (2 points)

   b)   Please give an expression for CAR(-1,1) under this regression model. (2 points)

   c)   Provide one advantage and one disadvantage of this regression model, compared to the for-loop version of the event study R code we studied in class. (2 points)