

A4 -Zihao Zhang

```
1. setwd("/Users/zhangzihao/Desktop/duke/613/HW/a4/Data")
2. library(data.table)
3. library(tidyverse)
4. dat_A4_panel <- fread("dat_A4_panel.csv")
5. dat_A4 <- fread("dat_A4.csv")
6.
7. # 1
8. # 1.1
9. dat_A4$age <- 2019 - dat_A4$KEY_BDATE_Y_1997
10. dat_A4$work_exp <- rowSums(dat_A4[,18:28],na.rm = "TRUE")
11. dat_A4$work_exp <- dat_A4$work_exp/52
12. dat_A4$work_exp <- round(dat_A4$work_exp,2)
```

	CV_WKSWK_JOB_DLI.11_2019	YSCH.3113_2019	YINC_1700_2019	age	work_exp	ed
	NA	NA	NA	38	0.00	
	NA	3	100000	37	12.42	
	NA	5	59000	36	1.69	
	NA	3	27000	38	1.92	
	NA	3	100000	37	13.46	
	NA	3	17000	37	2.25	
	NA	1	NA	36	2.37	
	NA	5	60000	38	4.19	
	NA	6	90000	37	3.23	

```
13.
14.
15. # 1.2
16. # I am a little bit confused about the question.
17. # Should the variable represent the total years of schooling for both residential p
    arents and biological parents?
18. # I Create three variables. total years of schooling of biological parents; total y
    ears of residential parents; total years of these four type persons.
19. dat_A4$CV_HGC_BIO_DAD_1997[which(dat_A4$CV_HGC_BIO_DAD_1997 == "95")] <- NA
20. dat_A4$CV_HGC_BIO_MOM_1997[which(dat_A4$CV_HGC_BIO_MOM_1997 == "95")] <- NA
21. dat_A4$CV_HGC_RES_DAD_1997[which(dat_A4$CV_HGC_RES_DAD_1997 == "95")] <- NA
22. dat_A4$CV_HGC_RES_MOM_1997[which(dat_A4$CV_HGC_RES_MOM_1997 == "95")] <- NA
23. dat_A4$edu_bio_parents <- rowSums(dat_A4[,8:9],na.rm = "TRUE")
24. dat_A4$edu_res_parents <- rowSums(dat_A4[,10:11],na.rm = "TRUE")
25. dat_A4$edu_parents <- rowSums(dat_A4[,8:11],na.rm = "TRUE")
```

26.

SWK_JOB_DLI.11_2019	YSCH.3113_2019	YINC_1700_2019	age	work_exp	edu_bio_parents	edu_res_parents	edu_parents
	NA	NA	38	0.00	24	24	48
	3	100000	37	12.42	32	29	61
	5	59000	36	1.69	12	12	24
	3	27000	38	1.92	24	12	36
	3	100000	37	13.46	24	24	48
	3	17000	37	2.25	12	12	24
	1	NA	36	2.37	12	12	24
	5	60000	38	4.19	18	18	36
	6	90000	37	3.23	18	18	36
	6	100000	35	5.08	18	18	36

27. # 1.3

28. library("ggplot2")

29. # 1.3.1

30. # firstly, i think i need to calculate the mean

31. # but then, i think box plot can be more intuitive

32. # income data by age groups

33. dat_plot <- dat_A4

34. income_age <- dat_A4 %>% group_by(age) %>% summarise(income = mean(YINC_1700_2019, na.rm = T))

35. p <- ggplot(data = income_age,

36. mapping = aes(

37. x = age,

38. y = income,

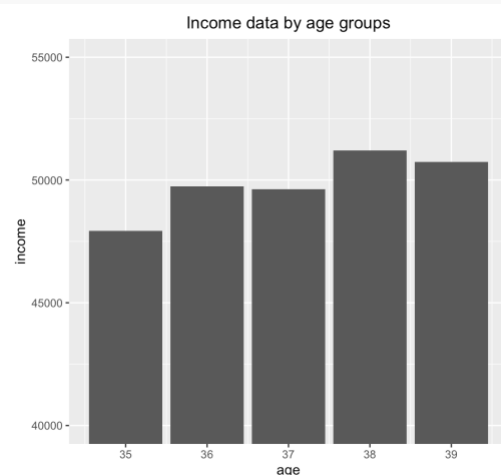
39.))

40. p + geom_col() +

41. coord_cartesian(ylim=c(40000,55000)) +

42. ggtitle("Income data by age groups") +

43. theme(plot.title = element_text(hjust = 0.5))



44.

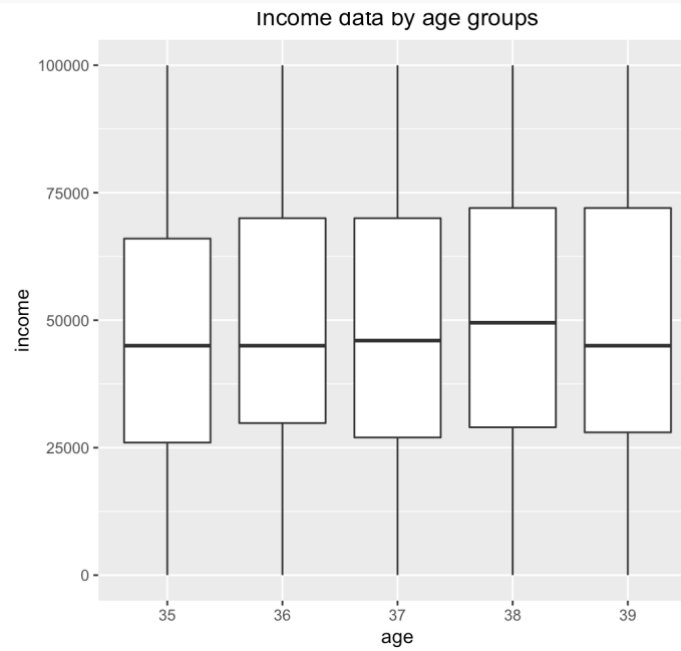
45. dat_A4 %>% dplyr::filter(!is.na(dat_A4\$YINC_1700_2019)) %>%

46. ggplot(aes(x = as.factor(age), y = YINC_1700_2019,)) + geom_boxplot() +

47. ggtitle("Income data by age groups") +

48. theme(plot.title = element_text(hjust = 0.5)) +

```
49. labs(x = "age", y = "income")
```



```
50.
```

```
51. # income data by gender groups
```

```
52. income_gender <- dat_A4 %>% group_by(KEY_SEX_1997) %>% summarise(income = mean(YINC_1700_2019, na.rm = T))
```

```
53. income_gender$KEY_SEX_1997 = c("male", "female")
```

```
54. p <- ggplot(data = income_gender,
```

```
55.       mapping = aes(
```

```
56.         x = KEY_SEX_1997,
```

```
57.         y = income
```

```
58.       )) +
```

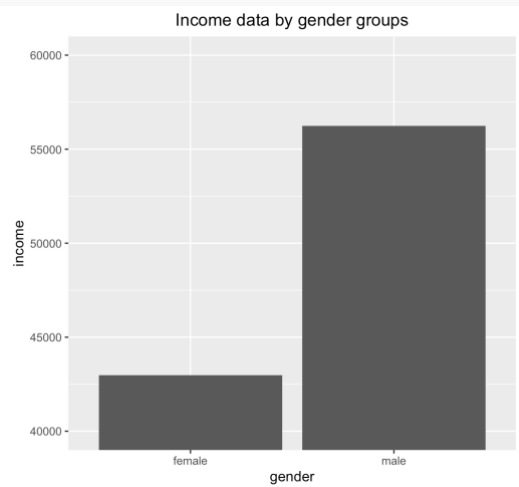
```
59.   labs(x = "gender")
```

```
60. p + geom_col() +
```

```
61.   coord_cartesian(ylim=c(40000, 60000)) +
```

```
62.   ggtitle("Income data by gender groups") +
```

```
63.   theme(plot.title = element_text(hjust = 0.5))
```



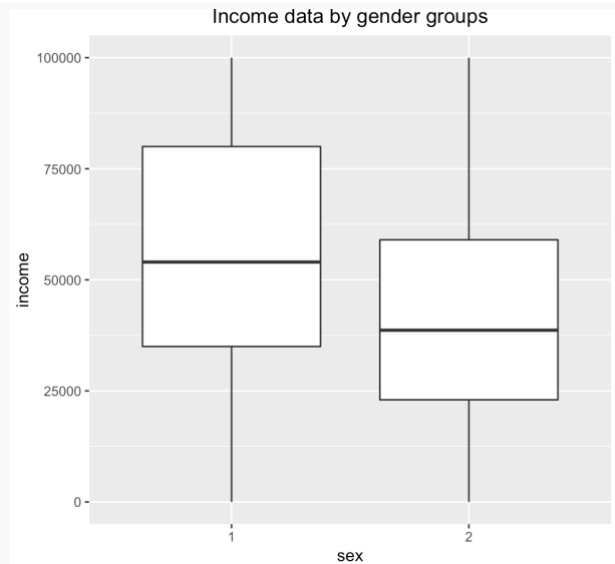
```
64.
```

```
65. dat_A4 %>% dplyr::filter(!is.na(dat_A4$YINC_1700_2019)) %>%
```

```

66. ggplot(aes(x = as.factor(KEY_SEX_1997), y = YINC_1700_2019 )) +
67.   geom_boxplot() +
68.   labs(x = "sex", y = "income") +
69.   ggtitle("Income data by gender groups") +
70.   theme(plot.title = element_text(hjust = 0.5))

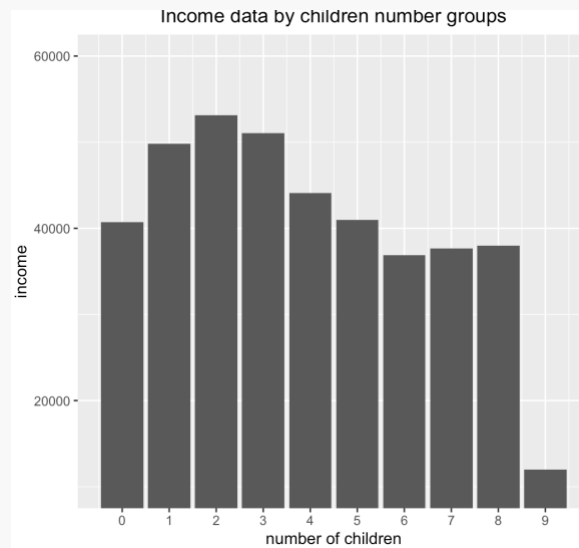
```



```

71.
72. # income data by number of children
73. income_child <- dat_A4 %>% group_by(CV_BIO_CHILD_HH_U18_2019) %>% summarise(income
  = mean(YINC_1700_2019,na.rm = T)) %>% na.omit()
74. p <- ggplot(data = income_child,
75.             mapping = aes(
76.               x = CV_BIO_CHILD_HH_U18_2019,
77.               y = income,
78.             ))+
79.   labs(x = "number of children")
80. p + geom_col() +
81.   coord_cartesian(ylim=c(10000,60000)) +
82.   scale_x_continuous(breaks = seq(0,9,1)) +
83.   ggtitle("Income data by children number groups") +
84.   theme(plot.title = element_text(hjust = 0.5))

```



85.

```
86. dat_A4 %>% dplyr::filter(!is.na(dat_A4$YINC_1700_2019)) %>%
```

```
87.   filter(CV_BIO_CHILD_HH_U18_2019 >= 0) %>%
```

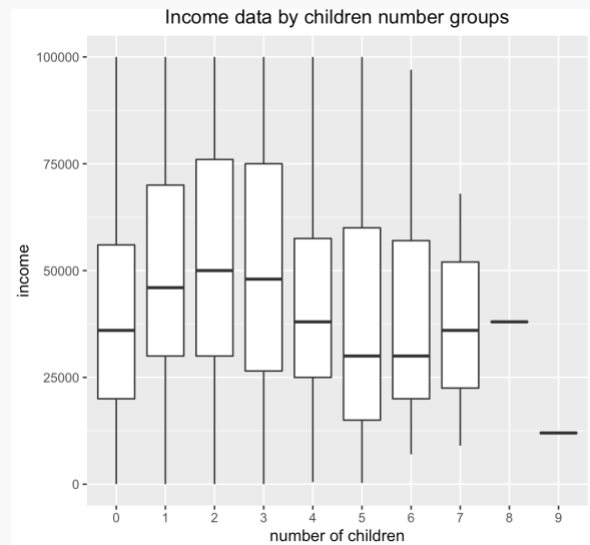
```
88.   ggplot(aes(x = as.factor(CV_BIO_CHILD_HH_U18_2019), y = YINC_1700_2019 )) +
```

```
89.   geom_boxplot() +
```

```
90.   labs(x = "number of children", y = "income") +
```

```
91.   ggtitle("Income data by children number groups") +
```

```
92.   theme(plot.title = element_text(hjust = 0.5))
```



93.

```
94. # 1.3.2
```

```
95. # share of "0" in the income data by age groups
```

```
96. share_inc_age <- dat_A4 %>% group_by(age) %>% summarize(share_zero = length(which((
  YINC_1700_2019==0))== 'TRUE'))/length(YINC_1700_2019))
```

```
97. p <- ggplot(data = share_inc_age,
```

```
98.   mapping = aes(
```

```
99.     x = age,
```

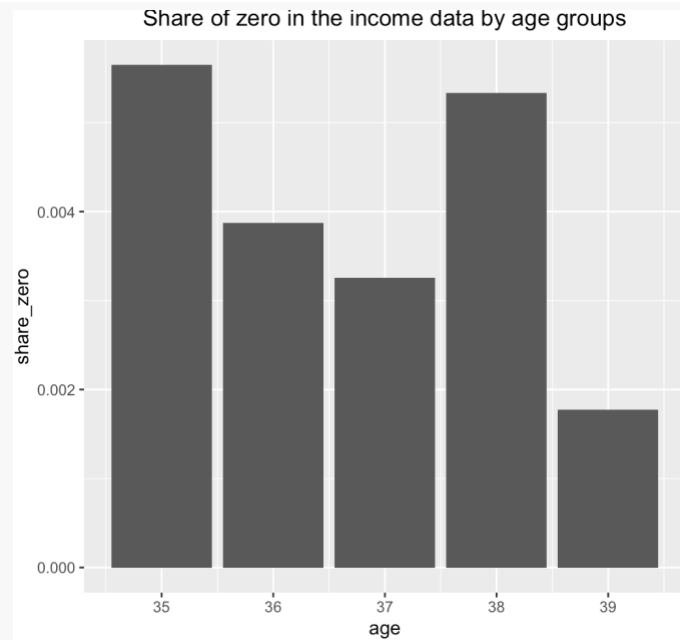
```
100.    y = share_zero,
```

```
101.  ))
```

```

102. p + geom_col() +
103.   ggtitle("Share of zero in the income data by age groups") +
104.   theme(plot.title = element_text(hjust = 0.5))

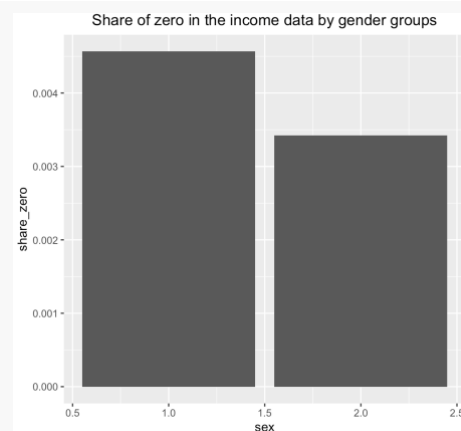
```



```

105.
106. # share of "0" in the income data by gender groups
107. share_inc_gender <- dat_A4 %>% group_by(KEY_SEX_1997) %>% summarize(share_zero = le
    ngth(which((YINC_1700_2019==0)=='TRUE'))/length(YINC_1700_2019))
108. p <- ggplot(data = share_inc_gender,
109.             mapping = aes(
110.               x = KEY_SEX_1997,
111.               y = share_zero,
112.             ))
113. p + geom_col() +
114.   ggtitle("Share of zero in the income data by gender groups") +
115.   theme(plot.title = element_text(hjust = 0.5)) +
116.   labs(x = "sex")

```



```

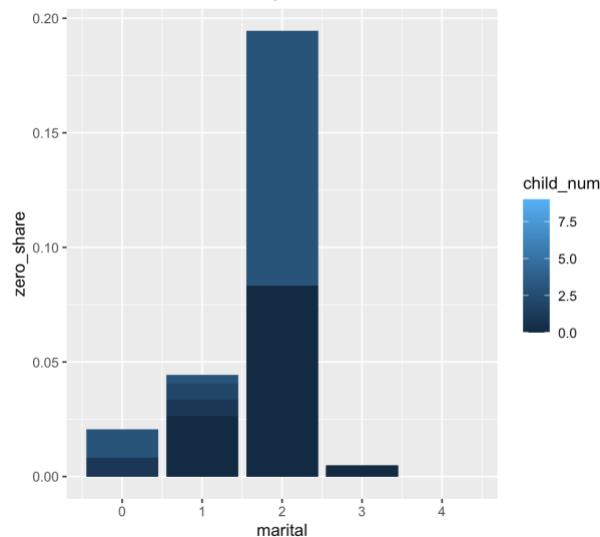
117.
118. # share of "0" in the income data by number of children and marital status

```

```

119. share_inc_child_marital <- dat_A4 %>% group_by(CV_MARSTAT_COLLAPSED_2019,CV_BIO_CHI
      LD_HH_U18_2019) %>%
120.   summarize(share_zero = length(which((YINC_1700_2019==0)=='TRUE'))/length(YINC_170
      0_2019))
121. colnames(share_inc_child_marital) <- c("marital", "child_num", "zero_share")
122. share_inc_child_marital <- na.omit(share_inc_child_marital)
123. p <- ggplot(data = share_inc_child_marital,
124.             mapping = aes(
125.               x = marital,
126.               y = zero_share,
127.               fill = child_num
128.             ))
129. p + geom_col() +
130.   ggtitle("Share of zero in the income data by children number and marital status")
131.   +
      theme(plot.title = element_text(hjust = 0.5))
      > of zero in the income data by children number and marital status

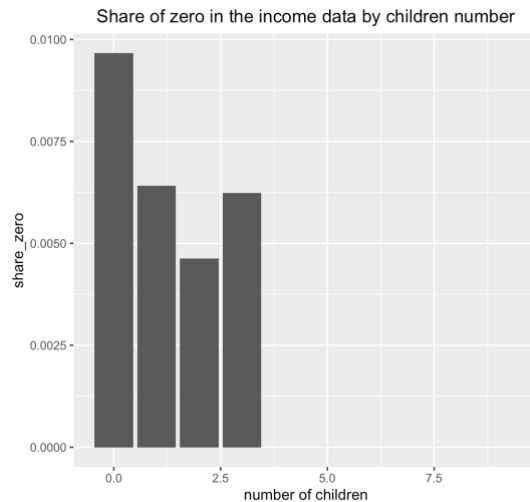
```



```

132.
133. # share of "0" in the income data by number of children
134. share_inc_child <- dat_A4 %>% group_by(CV_BIO_CHILD_HH_U18_2019) %>%
135.   summarize(share_zero = length(which((YINC_1700_2019==0)=='TRUE'))/length(YINC_170
      0_2019))
136. share_inc_child <- na.omit(share_inc_child)
137. p <- ggplot(data = share_inc_child,
138.             mapping = aes(
139.               x = CV_BIO_CHILD_HH_U18_2019,
140.               y = share_zero,
141.             ))
142. p + geom_col() +
143.   ggtitle("Share of zero in the income data by children number") +
144.   theme(plot.title = element_text(hjust = 0.5)) +
145.   labs(x = "number of children")

```



146.

```
147. # share of "0" in the income data by number of marital status
```

```
148. share_inc_marital <- dat_A4 %>% group_by(CV_MARSTAT_COLLAPSED_2019) %>%
```

```
149.   summarize(share_zero = length(which((YINC_1700_2019==0)=='TRUE'))/length(YINC_1700_2019))
```

```
150. share_inc_marital <- na.omit(share_inc_marital)
```

```
151. p <- ggplot(data = share_inc_marital,
```

```
152.   mapping = aes(
```

```
153.     x = CV_MARSTAT_COLLAPSED_2019,
```

```
154.     y = share_zero,
```

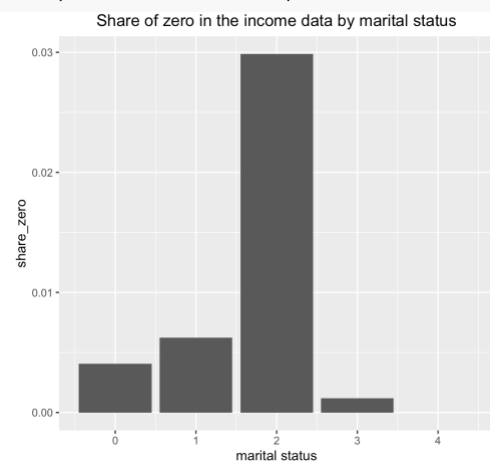
```
155.   ))
```

```
156. p + geom_col() +
```

```
157.   ggtitle("Share of zero in the income data by marital status") +
```

```
158.   theme(plot.title = element_text(hjust = 0.5)) +
```

```
159.   labs(x = "marital status")
```



160.

```
161. # 1.3.3
```

```
162. # Age and income show a positive correlation, but the correlation is not very significant
```

```
163. # Men earn significantly more than women
```

```
164. # The number of children is positively correlated with income at the beginning (with 1-3 children), and then becomes negative
```



```

165. # Families with 1-3 children seem to have the most income
166. # there is no obvious trend, it seems like that 35 and 38 years group has more zero
    proportion in income
167. # More men than women have no income
168. # Separated households have the largest share of no income
169. # More than half of these separated families have 2-3 children
170.
171. # 2
172. # 2.1
173. data_2.1 <- dat_A4 %>% filter(YINC_1700_2019 >0)
174.
175. # I set age/gender/work_exp/edu_parents as the independent variables
176. # using lm function to check
177. data_2.1 %>%
178.   lm(YINC_1700_2019 ~ work_exp + age + KEY_SEX_1997 + edu_parents, data =.) %>%
179.   summary()

```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	27229.95	9589.08	2.840	0.00453 **
work_exp	1092.08	66.38	16.453	< 2e-16 ***
age	462.10	256.58	1.801	0.07176 .
KEY_SEX_1997	-12747.51	713.41	-17.868	< 2e-16 ***
edu_parents	402.48	21.66	18.581	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 26120 on 5371 degrees of freedom
Multiple R-squared: 0.16, Adjusted R-squared: 0.1593
F-statistic: 255.7 on 4 and 5371 DF, p-value: < 2.2e-16

```

180.
181. # interpret the estimate results
182. # all independent variables are significant, The significance of "age" is relative
    y low, only significant at 95% significance
183. # work_exp: if other factors are fixed, workers earn more $1,092 for every addition
    al year of work experience
184. # age: if other factors are fixed, workers earn $462 more if one year older
185. # sex: if other factors are fixed, female earns $12747 less than male
186. # parents' education: if other factors are fixed, if parents' education years incre
    ase one year, the income increases $402
187.
188. # why there might be selection problem
189. # there are some interviewee reporting their incomes are zero or unwilling to repor
    t, which is may be not random.
190. # this phenomenon may influence the bias and the estimate of independent variables.
191.
192. # 2.2
193. # why Heckman can solve the problem

```

```

194. # To solve the selection problem, the Heckman model assumes that there are some other
      variables influencing dependent variable but not included in the independent variable
      sets.
195. # and the model estimates this part at the first stage
196. # then, using this part as a regressor at the second stage to avoid the problem
197.
198. # 2.3
199. data_2.3 <- dat_A4
200.
201. # create a dummy variables
202. for (i in 1:nrow(data_2.3)){
203.   if (isTRUE(data_2.3$YINC_1700_2019[i]>0)){
204.     data_2.3$dummy[i] = 1
205.   }else{
206.     data_2.3$dummy[i] = 0
207.   }
208. }
209. # missing value = 0
210.
211. # create inter
212. data_2.3$inter <- 1
213.
214. # create other variables
215. i <- data_2.3$inter
216. work_exp <- data_2.3$work_exp
217. edu <- data_2.3$edu_parents
218. age <- as.numeric(data_2.3$age)
219. sex <- as.numeric(data_2.3$KEY_SEX_1997)
220. d <- data_2.3$dummy
221.
222. # first stage
223. heck1 <- glm(formula = d ~ work_exp + age + edu + sex, family = binomial(link = "probit"),
                data = data_2.3)
224. summary(heck1)

```

```

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  0.3910549  0.4096167   0.955  0.33974
work_exp     0.2044746  0.0048510  42.151 < 2e-16 ***
age          -0.0297541  0.0109349  -2.721  0.00651 **
edu           0.0050751  0.0009089   5.584 2.35e-08 ***
sex           0.0001196  0.0305432   0.004  0.99687
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 12104.3  on 8983  degrees of freedom
Residual deviance: 8533.1  on 8979  degrees of freedom
AIC: 8543.1

```

```

226. predict1 <- -predict(heck1, newdata = NULL,
227.                        type = c("link", "response", "terms"),
228.                        se.fit = FALSE, dispersion = NULL, terms = NULL,
229.                        na.action = na.pass)
230.
231. imr <- (1/(1-pnorm(predict1))) * dnorm(predict1)
232.
233. heckfunc <- function (par, work_exp, edu, age, sex, imr) {
234.   yhat = par[1] + par[2]* work_exp + par[3]* edu + par[4] * age + par[5] * sex + p
         ar[6] * imr
235.   prob = pnorm(yhat)
236.   prob[prob>0.999999] = 0.999999
237.   prob[prob<0.000001] = 0.000001
238.   like = imr*log(prob) + (1-imr)*log(1-prob)
239.   return( - sum(like) )
240. }
241.
242. datah <- cbind(data_2.3,imr)
243. imr_reg = datah$imr
244.
245. predictor <- lm(YINC_1700_2019 ~ work_exp + edu + age + sex + imr, data = datah)
246. summary(predictor)

```

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  35910.81    9545.92   3.762  0.00017 ***
work_exp     -619.28     150.52  -4.114  3.94e-05 ***
edu           321.96      22.27  14.458 < 2e-16 ***
age          1006.09     257.44   3.908  9.42e-05 ***
sex        -12796.75     707.62 -18.084 < 2e-16 ***
imr         -30100.47    2369.76 -12.702 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 25990 on 5406 degrees of freedom
(3572 observations deleted due to missingness)
Multiple R-squared:  0.18,    Adjusted R-squared:  0.1792
F-statistic: 237.3 on 5 and 5406 DF,  p-value: < 2.2e-16

```

```

247.
248. start <- predictor$coefficients
249. results <- optim(start, fn = heckfunc, method = "BFGS",
250.                  control = list(trace = 6, maxit = 3000),
251.                  work_exp = work_exp, edu = edu, age = age, sex = sex, imr =
                     imr_reg)
252. results$par
253.

```

```

initial value 39410.827908
final value 39410.827908
converged
> results$par
(Intercept)  work_exp      edu      age      sex      imr
35910.8089   -619.2836   321.9631  1006.0943 -12796.7481 -30100.4700

```

```
254. # The difference is obvious, so we can assume that the absence of the data of low e
      arning or zero may be not random.
```

```
255. # It influences the result of OLS
```

```
256. # the coefficients of work experience and sex are negative showing a negative
      correlation between these variables and income. And the edu and age have a positive
      correlation with income. Specifically, for example, when all other conditions
      fixed, male earns $12796 more than female.
```

```
257.
```

```
258. # 3
```

```
259. # 3.1
```

```
260. p <- ggplot(data = dat_A4,
```

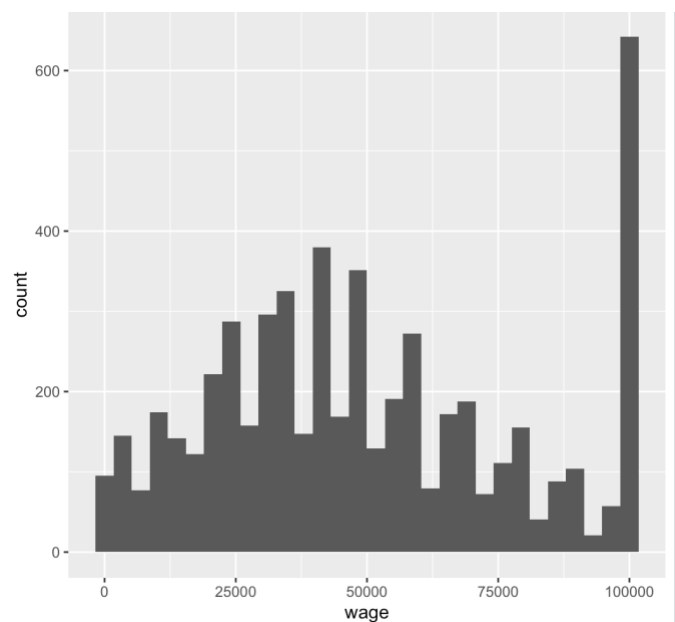
```
261.           mapping = aes(
```

```
262.             x = YINC_1700_2019,
```

```
263.           ))
```

```
264. p + geom_histogram() +
```

```
265.   labs(x = "wage")
```



```
266.
```

```
267. # what might be the censored value
```

```
268. # the income over 100000 should be top-coded
```

```
269.
```

```
270. # 3.2
```

```
271. # we can use a two-step model to solve the censoring problem
```

```
272. # Similar with question 2.3, we can firstly explain top-
      coded incidents and then we can use the inverse mills ratio to estimate at the seco
      nd stage.
```

```
273. # 3.3
```

```
274. data_3.2 <- dat_A4
```

```
275.
```

```
276. for (i in 1:nrow(data_3.2)){
```

```

277. if (isTRUE(data_3.2$YINC_1700_2019[i]<100000)){
278.   data_3.2$d[i] = 1
279. }else{
280.   data_3.2$d[i] = 0
281. }
282. }
283.
284. data_3.2$inter <- 1
285. inter <- data_3.2$inter
286. d2 <- data_3.2$d
287. work_exp <- data_3.2$work_exp
288. edu <- data_3.2$edu_parents
289. age <- as.numeric(data_3.2$age)
290. sex <- as.numeric(data_3.2$KEY_SEX_1997)
291. inc <- data_3.2$YINC_1700_2019
292.
293. # we may use Tobit model to solve the censoring problem
294. library(AER)
295. tobit <- tobit(inc ~ work_exp + edu + age + sex, right = 100000, data = data_3.2)
296. summary(tobit)
297.

```

```

Coefficients:
      Estimate Std. Error z value Pr(>|z|)
(Intercept)  2.053e+04  1.093e+04   1.878   0.0604 .
work_exp      1.184e+03  7.572e+01  15.633 <2e-16 ***
edu           4.481e+02  2.477e+01  18.092 <2e-16 ***
age           6.571e+02  2.926e+02   2.246   0.0247 *
sex          -1.401e+04  8.134e+02 -17.230 <2e-16 ***
Log(scale)    1.029e+01  1.066e-02  966.074 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Scale: 29576

Gaussian distribution
Number of Newton-Raphson Iterations: 3
Log-likelihood: -5.633e+04 on 6 Df
Wald-statistic: 943.5 on 4 Df, p-value: < 2.22e-16

```

```

298. par <- tobit$coefficients
299. par[6] <- 10.295
300. par <- as.vector(c(par))
301.
302. tlikefunc <- function(par, work_exp, edu, age, sex, d2, inc){
303.   yhat = par[1] + par[2]* work_exp + par[3]* edu + par[4] * age + par[5] * sex
304.   like = (1-d2) * log(1 - pnorm((100000 - yhat)/exp(par[7]))) + d2*log(dnorm((inc-
      yhat)/exp(par[7]))/exp(par[7]))
305.   return(-sum(like,na.rm=T))
306. }
307.
308. start = par + runif(6,-1,1)

```

```

309.results <- optim(start, fn = tlikefunc, method = "BFGS",
310.                  control = list(trace = 6, maxit = 3000),
311.                  work_exp = work_exp, edu = edu, age = age, sex = sex, d2 = d2, inc
    = inc)

```

```
312.results$par
```

```
313.
```

```

initial value -0.000000
final value -0.000000
converged
> results$par
[1] 20524.76250 1183.41274 449.09457 657.10549 -14015.17576 10.83978

```

```
314.# not correcting for the censored data
```

```
315.ols <- lm(inc ~ work_exp + edu + age +sex, data = data_3.2)
```

```
316.summary(ols)
```

```
317.ols$coefficients
```

```
318.results$par
```

```
319.
```

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  25150.3      9648.2   2.607  0.00917 **
work_exp      1100.0        66.8  16.468 < 2e-16 ***
edu           396.4         21.8  18.186 < 2e-16 ***
age           508.3        258.2   1.969  0.04904 *
sex          -12590.9       717.9 -17.540 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 26370 on 5407 degrees of freedom
(3572 observations deleted due to missingness)
Multiple R-squared:  0.1555,    Adjusted R-squared:  0.1549
F-statistic: 248.9 on 4 and 5407 DF,  p-value: < 2.2e-16

> ols$coefficients
(Intercept)  work_exp      edu      age      sex
25150.3154  1100.0225  396.4440  508.2759 -12590.8694

```

```
320.# 3.4
```

```

321.# we can find the absolute value of the coefficient of each variable is smaller if
    we ignore the censoring problem

```

```
322.# it means that we underestimate the effect.
```

```
323.
```

```
324.# 4
```

```
325.# 4.1
```

```
326.# Correlations between variables can lead to selection problems.
```

```

327.# For example, people with better educational background may have better family bac
    kground and talents, which makes them more popular in the marriage market and easie
    r to obtain better income.

```

```

328.# In addition, good marital status may allow people to focus more on work, leading
    to higher earnings.

```

```
329.# These situations can lead to the selection problem.
```

```
330.
```

```
331.# 4.2
```

```
332.# data preparing
```

```

333. data4 <- dat_A4_panel
334.
335. # it is hard to conduct based on initial data
336. # so i need to convert the wide data to long
337. # firstly i need to rename the variables
338. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_1998"] <-
      "CV_HIGHEST_DEGREE_9899_1998"
339. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_1999"] <-
      "CV_HIGHEST_DEGREE_9900_1999"
340. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2000"] <-
      "CV_HIGHEST_DEGREE_0001_2000"
341. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2001"] <-
      "CV_HIGHEST_DEGREE_0102_2001"
342. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2002"] <-
      "CV_HIGHEST_DEGREE_0203_2002"
343. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2003"] <-
      "CV_HIGHEST_DEGREE_0304_2003"
344. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2004"] <-
      "CV_HIGHEST_DEGREE_0405_2004"
345. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2005"] <-
      "CV_HIGHEST_DEGREE_0506_2005"
346. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2006"] <-
      "CV_HIGHEST_DEGREE_0607_2006"
347. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2007"] <-
      "CV_HIGHEST_DEGREE_0708_2007"
348. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2008"] <-
      "CV_HIGHEST_DEGREE_0809_2008"
349. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2009"] <-
      "CV_HIGHEST_DEGREE_0910_2009"
350. names(data4)[names(data4) == "CV_HIGHEST_DEGREE_EVER_EDT_2010"] <-
      "CV_HIGHEST_DEGREE_1011_2010"
351.
352. library(panelr)
353. data4 <- long_panel(
354.   data4,
355.   prefix = '_',
356.   suffix = NULL,
357.   begin = 1997,
358.   end = 2019,
359.   id = "id",
360.   wave = "wave",
361.   periods = NULL,
362.   label_location = "end",
363.   as_panel_data = TRUE,

```

```

364. match = ".*",
365. use.regex = FALSE,
366. check.varying = TRUE
367. )
368.
369. data4 <- data4 %>%
370.   filter(wave!='2012') %>%
371.   filter(wave!='2014') %>%
372.   filter(wave!='2016') %>%
373.   filter(wave!='2018')
374.
375. # new variables
376. # work_experience
377. we <- data4[,c(10:16,20:27)]
378. we <- as.matrix(we)
379. we[is.na(we)] <- 0
380. we <- we[,3:17]
381. we <- as.data.frame(we)
382. for (i in 1:14) {
383.   we[,i]<-as.numeric(we[,i])
384. }
385.
386. data4$work_exp <- rowSums(we[,1:14])/52
387. data4$work_exp <- round(data4$work_exp,2)
388.
389. # age
390. data4$age <- data4$wave - data4$KEY_BDATE_Y
391.
392. # between estimator: gender, work_exp, edu, marital status
393. m_gender <- data4 %>% group_by(id) %>% summarize(m_gender=mean(KEY_SEX,na.rm = TRUE
  ))
394. m_work_exp <- data4 %>% group_by(id) %>% summarize(m_work_exp=mean(work_exp,na.rm
  = TRUE))
395. m_ms <- data4 %>% group_by(id) %>% summarize(m_ms=mean(CV_MARSTAT_COLLAPSED,na.rm
  = TRUE))
396.
397. data4$'YINC-1700' <- as.numeric(data4$'YINC-1700')
398. m_inc <- data4 %>% group_by(id) %>% summarize(m_inc=mean(`YINC-
  1700`,na.rm = TRUE))
399. data4$CV_HIGHEST_DEGREE_EVER_EDT <- as.numeric(data4$CV_HIGHEST_DEGREE_EVER_EDT)
400. m_edu <- data4 %>% group_by(id) %>% summarize(m_edu=mean(CV_HIGHEST_DEGREE_EVER_EDT
  ,na.rm = TRUE))
401.
402. between <- lm(m_inc$m_inc~ m_edu$m_edu + m_work_exp$m_work_exp + m_ms$m_ms)

```



```
403. summary(between)
```

```
404.
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    7647.95     387.56  19.734 < 2e-16 ***
m_edu$m_edu     3446.72     115.91  29.735 < 2e-16 ***
m_work_exp$m_work_exp 2302.74      96.63  23.830 < 2e-16 ***
m_ms$m_ms       1958.12     335.18   5.842 5.36e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 14510 on 7998 degrees of freedom
(982 observations deleted due to missingness)
Multiple R-squared:  0.1943,    Adjusted R-squared:  0.194
F-statistic: 643.1 on 3 and 7998 DF,  p-value: < 2.2e-16
```

```
405. # Within Estimator: exper, edu, marital status
```

```
406. data4 <- left_join(data4,m_inc,by = "id")
```

```
407. data4 <- left_join(data4,m_edu,by='id')
```

```
408. data4 <- left_join(data4,m_work_exp,by='id')
```

```
409. data4 <- left_join(data4,m_ms,by='id')
```

```
410.
```

```
411. data4$d_inc <- data4$YINC-1700'-data4$m_inc
```

```
412. data4$d_edu <- data4$CV_HIGHEST_DEGREE_EVER_EDT - data4$m_edu
```

```
413. data4$d_work_exp <- data4$work_exp - data4$m_work_exp
```

```
414. data4$d_marital_status <- data4$CV_MARSTAT_COLLAPSED - data4$m_ms
```

```
415.
```

```
416. within <- lm(d_inc ~ d_edu + d_work_exp + d_marital_status, data = data4)
```

```
417. summary(within)
```

```
418.
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   12749.84     256.98   49.61 <2e-16 ***
d_edu          9815.15     674.39   14.55 <2e-16 ***
d_work_exp     1405.78      49.28   28.53 <2e-16 ***
d_marital_status 4108.70     302.45   13.59 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 27990 on 25813 degrees of freedom
(144879 observations deleted due to missingness)
Multiple R-squared:  0.04942,    Adjusted R-squared:  0.04931
F-statistic: 447.3 on 3 and 25813 DF,  p-value: < 2.2e-16
```

```
419. # first difference estimator
```

```
420. dataf <- data4
```

```
421. dataf$first_inc <- ave(dataf$YINC-1700', dataf$id, FUN=function(x)dplyr::lag(x))
```

```
422. dataf$first_work_exp <- ave(dataf$work_exp, dataf$id, FUN=function(x)dplyr::lag(x))
```

```
423. dataf$first_edu <- ave(dataf$CV_HIGHEST_DEGREE_EVER_EDT, dataf$id, FUN=function(x)dplyr::lag(x))
```

```
424. dataf$first_ms <- ave(dataf$CV_MARSTAT_COLLAPSED, dataf$id, FUN=function(x)dplyr::lag(x))
```

```
425.
```

```

426. dataf$fd_inc <- dataf$'YINC-1700' - dataf$first_inc
427. dataf$fd_edu <- dataf$CV_HIGHEST_DEGREE_EVER_EDT - dataf$first_edu
428. dataf$fd_work_exp <- dataf$work_exp - dataf$first_work_exp
429. dataf$fd_ms <- dataf$CV_MARSTAT_COLLAPSED - dataf$first_ms
430.
431. fd <- lm(fd_inc ~ fd_work_exp + fd_edu + fd_ms, data = dataf)
432. summary(fd)
433.

```

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  6190.28    207.64   29.813 < 2e-16 ***
fd_work_exp   459.25     69.77    6.582 4.77e-11 ***
fd_edu       -486.66    698.79   -0.696  0.486
fd_ms         232.11    419.11    0.554  0.580
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 24660 on 16387 degrees of freedom
(154305 observations deleted due to missingness)
Multiple R-squared:  0.002681, Adjusted R-squared:  0.002498
F-statistic: 14.68 on 3 and 16387 DF, p-value: 1.517e-09

```

```

434. # 4.3
435. # all variables in between and within estimator are significant and positive.
436. # the results show that when other condition fixed, the increase of the increase le
      ading to increase of income
437. # for example, the coefficient of work experience in between estimator is 1958, whi
      ch means, when all conditions fixed, if the mean of work experience increase one ye
      ar, the mean of income increase $1958
438. # two of variables in first estimator are not significant, it confused me.
439. # the difference may be caused by different method of dealing with NA
440. # This difference is more pronounced in the process of the first difference estimat
      or, which may be the reason why this model is not significant

```