

Assignment 1

Walter Verwer & Bas Machielsen

January 5, 2021

Question 1: The sample selection model.

A researcher aims to gain insight in the potential earnings of the non-employed. (In the data, the non-employed can be identified by a missing value for the earnings variable). She realizes that the sample of observed wages may be subject to sample selection.

(a) Run an OLS regression for log-earnings on schooling, age, and age squared. Present the results and comment on the estimates.

```
model_1 <- lm(data = data, formula = logWage ~ schooling + age + age2)

results_1 <- summary.lm(model_1)
coeffs_1 <- results_1$coefficients[,1]
pvals_1 <- results_1$coefficients[,4]

stargazer(model_1, style = "AER",
           font.size = "small",
           header = F, label = 'tab:q1_a_ols',
           title = 'OLS regression for log-earnings on schooling, age and age squared.')
```

The results show that one additional year of schooling has an effect of 0.216 on $\log(\text{Wage})$, which means that one additional year of schooling has an estimated 21.6% effect on wages earned. This result is highly significant (p-value is 2.706×10^{-11}). Another result, shown in figure 1 is that being one year older has an estimated effect of -0.342 on $\log(\text{Wage})$. Thus, being one year older is estimated to have a -34.189% on wages. This result is not significant at a common level (p-value is 0.512). The third estimated coefficients is the one corresponding to the variable age squared. It has an estimated coefficient of -0.011, which represents the estimated effect of a one unit increase in age squared on $\log(\text{Wage})$. This also means that a one unit increase in age squared is equal to an estimate effect of -1.114% on the linear representation of wage. This effect is not significant at common levels, because the p-value is 0.184. Finally, the estimated constant in the model is estimated at 26.409. This means that if all other variables are zero, then the $\log(\text{Wage})$ will be equal to 26.409. Thus, if all other variables are zero, then wage is estimated to be equal to $\exp(26.409) = 2.947 \times 10^{11}$. This is an extremely high number given the characteristics of the wage variable. However, this result is highly significant, because the p-value is 0.001.

(b) Briefly discuss the sample selection problem that may arise in using these OLS estimates for the purpose of predicting the potential earnings of the non-employed. Formulate the sample selection model. In your answer, include an explanation why OLS may fail in this context.

An individual is only in this dataset if they earn wages, i.e. if they are employed. Being employed itself is not randomly allocated, but rather, a function of e.g. age, age², and schooling. Also, employment status follows from labor supply and demand forces. Hence, the estimates of schooling on earnings are conditional on having earnings to begin with, whereas unbiased estimates must also include those individuals.

Table 1: OLS regression for log-earnings on schooling, age and age squared.

	logWage
schooling	0.216*** (0.032)
age	-0.342 (0.521)
age2	-0.011 (0.008)
Constant	26.400*** (8.060)
Observations	416
R ²	0.815
Adjusted R ²	0.813
Residual Std. Error	1.500 (df = 412)
F Statistic	604.000*** (df = 3; 412)
Notes:	***Significant at the 1 percent level. **Significant at the 5 percent level. *Significant at the 10 percent level.

Formally, $\mathbb{E}[\text{Earnings}] = \mathbb{E}[\text{Earnings}|\text{Having a job}] \cdot \mathbb{P}[\text{Having a job}] + \mathbb{E}[\text{Earnings}|\text{Not having a job}] \cdot (1 - \mathbb{P}[\text{Having a job}])$. The given estimation only concerns $\mathbb{E}[\text{Earnings}|\text{Having a job}]$.

The sample selection model is given by the following two equations.

$$\mathbb{P}[I_i^* = 1|Z] = \Phi(Z_i'\gamma) \quad (1)$$

In equation 1, I_i^* denotes an indicator variable which is equal to 1 if we observe the wage of an individual, $\Phi(\cdot)$ denotes a normal CDF, and Z_i denotes a vector of explanatory variables for the probability of an individual being employed or not. Note that the model used in this situation concerns a probit model, which has the goal to estimate the probability that an individual is employed.

The second equation is concerned with explaining the wage of an individual. It is given by the following equation.

$$Y_i^* = X_i'\beta + U_i \quad (2)$$

In equation 2, Y_i^* denotes the observed $\log(\text{Wage})$, X_i are the explanatory regressors for explaining $\log(\text{Wage})$.

(c) Which variable in your data may be a suitable candidate as an exclusion restriction for the sample selection model? For an exclusion restriction variable, we need a variable that is significantly away from zero in the selection equation and does not have an effect in the equation intended to explain the potential earnings of the non-employed. A potential candidate is this a variable that matters for being employed or not, but does not influence the height of the wage. We believe that marriage status could be a suitable candidate variable. The reason being that marriage status could matter for being employed or not, because if one is not married, the person is more likely to be the only person that has to provide the necessary funds of living. If someone is married however, than it is more likely that this person is not employed, because it might be the case that the partner provides. Marriage status is arguably a variable that does not provide a direct effect on the height of wages earned. Seeing as both criteria are met, we conclude that marriage status is a suitable candidate for an exclusion restriction.

(d) Estimate the sample selection model with the Heckman two-step estimator, both with and without the exclusion restriction and compare the outcomes.

```
# Construct I (I=1 for y_i* != na, else 0)
data$I <- ifelse(is.na(data$logWage) , 0, 1) # if na, then I is zero. Else 1.

## Two stage approach:
# Stage 1:
probit <- glm(I ~ married,
              family = binomial(link = "probit"),
              data = data)

# Construct variable for stage 2:
Z_gamma <- cbind(1,data$married) %*% coef(probit)
inverse_mills_ratio <- dnorm(Z_gamma, mean=0, sd=1) / pnorm(Z_gamma, mean=0, sd=1)

# Stage 2:
sample_selection_2sls <- lm(logWage ~ schooling + age + age2 + inverse_mills_ratio,
                           data=data)

## using package 'sampleSelection' to get maximum likelihood estimates:
sample_selection_ml <- heckit(I ~ married,
                             logWage ~ schooling + age + age2, data=data, method='ml')

# Obtain output:
stargazer(sample_selection_2sls,
          se = starprep(sample_selection_2sls, se_type = "HC1"),
          style = "AER",
          font.size = "small",
          header = F, label = 'tab:q1_d_2sls',
          title = 'Sample selection regression for log-earnings on schooling, age
and age squared, using two stage approach and maximum likelihood')
```

(e) Estimate the sample selection model with Maximum Likelihood, both with and without the exclusion restriction and compare the outcomes.

(f) On the basis of your results, how would you specify the distribution of potential earnings for the non-employed?

Question 2: Earnings and Schooling

The same researcher is interested in estimating the causal effect of schooling on earnings for employed individuals only. As a consequence, she performs the subsequent analysis on the (sub)sample of employed individuals.

(a) Discuss the estimation of the causal effect of schooling on earnings by OLS. In particular, address whether or not it is plausible that regularity conditions for applying OLS are satisfied.

It is not plausible that the regularity conditions are satisfied. In particular, an observable such as an individual's **ability** might be correlated with the wage, but also with the decision to live close to a school. Hence, the estimates suffer from endogeneity.

(b) The researcher has collected data on two potential instrumental variables subsidy and distance for years of schooling.

Table 2: Sample selection regression for log-earnings on schooling, age and age squared, using two stage approach and maximum likelihood

	logWage
schooling	0.216*** (0.031)
age	-0.352 (0.475)
age2	-0.011 (0.008)
inverse_mills_ratio	-0.150 (0.609)
Constant	26.600*** (7.470)
Observations	416
R ²	0.815
Adjusted R ²	0.813
Residual Std. Error	1.500 (df = 411)
F Statistic	452.000*** (df = 4; 411)
Notes:	***Significant at the 1 percent level. **Significant at the 5 percent level. *Significant at the 10 percent level.

- distance measures the distance between the school location and the residence of the individual while at school-going age.
- subsidy is an indicator depending on regional subsidies of families for covering school expenses.

The researcher has the option to use only distance as an instrumental variable, or to use only the instrumental variable subsidy, or to use both distance and subsidy as instrumental variables. Perform instrumental variables estimation for these three options. Which option do you prefer? Include in your answer the necessary analyses and numbers on which you base your choice.

```
firstoption <- ivreg(data = data, formula =
  logWage ~ age + age2 + schooling | distance + age + age2)

secondoption <- ivreg(data = data, formula =
  logWage ~ age + age2 + schooling | subsidy + age + age2)

thirdoption <- ivreg(data = data, formula =
  logWage ~ age + age2 + schooling | subsidy + distance + age
  + age2)

stargazer(firstoption, secondoption, thirdoption, font.size = "small",
  style = "AER",
  header = F)
```

We consider that the second option, to include only subsidy as an instrument, is the best option. The reason is that distance is unlikely to satisfy the exclusion restriction: distance is (to a certain extent) an

Table 3:

	logWage		
	(1)	(2)	(3)
age	−0.192 (0.587)	−0.233 (0.546)	−0.229 (0.547)
age2	−0.014 (0.010)	−0.013 (0.009)	−0.013 (0.009)
schooling	0.470 (0.299)	0.401*** (0.106)	0.408*** (0.102)
Constant	22.700** (9.700)	23.700*** (8.520)	23.600*** (8.530)
Observations	416	416	416
R ²	0.786	0.799	0.798
Adjusted R ²	0.784	0.798	0.797
Residual Std. Error (df = 412)	1.610	1.560	1.560
Notes:	***Significant at the 1 percent level. **Significant at the 5 percent level. *Significant at the 10 percent level.		

endogenous variable: wealthier (or more able) parents may choose to live closer to school, and invest more in the education of their children (or genetically transmit ability). Since a potentially endogenous instrument must not be used as such, we prefer the estimates in equation 2. However, we see that the results show that distance has no predictive power in schooling, thus showing that the endogeneity is very small. Conditional on subsidy being a good instrument, then, the potential endogeneity does not substantially changes the estimates of schooling on earnings.

- (c) Compare the IV estimates with the OLS outcomes. Under which conditions would you prefer OLS over IV? Perform a test and use the outcome of the test to support your choice between OLS and IV. Motivate your choice.

We first observe that the OLS estimate $\beta = 0.216$ is about half the magnitude of the IV-estimate. This means that the bias generated by OLS likely *downplays* the actual effect (if the IV estimates satisfy the exclusion restriction). In case we would not trust the IV assumptions, we would prefer to trust the (conservative) estimate that downplays the effect, i.e. the OLS estimates. We can test whether the OLS estimates are substantially different from the IV estimates by conducting a Hausman test:

```

hoi <- summary(
  secondoption,
  diagnostics=TRUE)

hoi$diagnostics

```

```

##              df1 df2 statistic  p-value
## Weak instruments    1 412     43.32 1.42e-10
## Wu-Hausman         1 411      3.63 5.73e-02
## Sargan              0  NA         NA      NA

```

The null hypothesis in the Hausman test is exogeneity of the *schooling* variable. As becomes clear, the null hypothesis is marginally rejected, implying the *schooling* is endogenous, but only marginally so. Hence, we would prefer to trust the IV estimates in this case.