

# Assignment 1

Students

January 6, 2023

## Question 1

i)

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

	<i>Dependent variable:</i>
	logwage
schooling	0.216*** (0.032)
age	-0.342 (0.521)
I(age^2)	-0.011 (0.008)
Constant	26.409*** (8.057)
Observations	416
R <sup>2</sup>	0.815
Adjusted R <sup>2</sup>	0.813
Residual Std. Error	1.499 (df = 412)
F Statistic	604.261*** (df = 3; 412)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

From Table 1 can be observed that only the intercept and **schooling** are significant. Both are significant at the 1%-significance level. For a given worker, an additional year of schooling is associated with a  $(e^{0.216} - 1) \cdot 100 \approx 24.11\%$  increase in wage. The intercept and explanatory variables explain 81.5% of the variation in **logwage**.

ii)

Let  $Z'_i$ ,  $\gamma$ , and  $V_i$  denote one of the regressors **schooling**, **age** and **I(age^2)** or an exclusion restriction (in case of correct specification), the corresponding coefficient, and error term for individual  $i$ , respectively. Then, we have the selection equation

$$I_i^* = Z'_i \gamma + V_i,$$

where  $I_i$  takes value 1 if the wage for individual  $i$   $I_i^* > 0$ , and value 0 otherwise. In addition, let  $Y_i^*$  denote the latent variable **logwage**,  $X_i'$  the regressors,  $\beta$  the coefficient vector and  $U_i$  the error term. The second regression equation is

$$Y_i^* = X_i' \beta + U_i.$$

However, instead of the true **logwage**  $Y_i^*$  we observe

$$Y_i = \begin{cases} Y_i^* & \text{if } I_i = 1 \\ \text{missing} & \text{if } I_i = 0. \end{cases}$$

Hence, **logwage** is only observed when a worker earns a wage greater than 0. In the 2-dimensional case (1 regressor), observations are essentially missing at the bottom of the y-axis, which will result in biased OLS estimates when running a regression. Obviously, this problem extends to higher dimensions. Since the problem is not alleviated when obtaining a larger sample, the OLS estimates will also be inconsistent.

### iii)

The exclusion restriction must explain whether observations are in the sample. That is, whether an individual earns a wage. However, the exclusion restriction should not explain variation in **logwage**. Therefore, it can be argued that **married** is a suitable candidate. Whether an individual is married will not directly affect the (potential) wage they may earn. However, it can affect whether an individual is in the labor force. For example, married individuals may choose not to work to take care of their families.

Table 2: Two-step Heckman estimation without exclusion restriction.

	Dependent variable:
	logwage
schooling	0.915
age	-0.347 (0.592)
I(age^2)	-0.011 (0.010)
Constant	-16.019
Observations	666
$\rho$	1.334
Inverse Mills Ratio	62.404
Note:	*p<0.1; **p<0.05; ***p<0.01

Table 2 shows the Two-step Heckman estimation where the selection variable  $Z'$  is **schooling**, whereas Table 3 shows the Two-step Heckman estimation where the selection variable  $Z'$  is **married**. In the first case, where  $Z'$  is a variable in  $X'$ , there is perfect collinearity (due to **schooling**) in the second step of the estimation. As a consequence, the standard error of **schooling** is **Inf**, and is not reported. In the second step, where an appropriate exclusion restriction is used, both the coefficients and the standard errors are obtained.

Table 3: Two-step Heckman estimation with exclusion restriction.

	<i>Dependent variable:</i>
	logwage
schooling	0.216*** (0.031)
age	−0.352 (0.520)
I(age <sup>2</sup> )	−0.011 (0.008)
Constant	26.648*** (8.075)
Observations	666
$\rho$	−0.100
Inverse Mills Ratio	−0.150 (0.607)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

iv)

Table 4: Maximum likelihood estimation without exclusion restriction.

	<i>Dependent variable:</i>
	logwage
schooling	0.061 (Inf.000)
age	1.436 (Inf.000)
I(age <sup>2</sup> )	−0.037 (Inf.000)
Constant	−15.898 (Inf.000)
Observations	666
Log Likelihood	−1,732.567
$\rho$	1.000 (Inf.000)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Table 4 shows the maximum likelihood estimation where the selection variable  $Z'$  is `schooling`, whereas Table 5 shows the maximum likelihood estimation where the selection variable  $Z'$  is `married`. From the tables, we draw

Table 5: Maximum likelihood estimation with exclusion restriction.

	<i>Dependent variable:</i>
	logwage
schooling	0.216*** (0.031)
age	−0.350 (0.520)
I(age^2)	−0.011 (0.008)
Constant	26.616*** (8.063)
Observations	666
Log Likelihood	−1,187.699
$\rho$	−0.089 (0.376)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

the same conclusion as in iii). However, for the maximum likelihood estimation all standard errors are reported to be `Inf` when the selection variable  $Z'$  is `schooling`, instead of only the standard error of `schooling`.

v)

## Question 2

i)

An OLS model of the form

$$earnings = \beta_0 + \beta_1 \text{schooling} + \epsilon$$

likely does not yield consistent estimates as the exogeneity assumption  $E[\mathbf{X}\epsilon]$  is going to be violated. In this case, the residual contains individual's characteristics such as intelligence, which should both influence the amount of schooling an individual seeks, as well as directly impact the income of an individual.

ii)

Table 6 displays the results for the different IV models and the OLS estimation.

Based on the F-statistics of the IV regressions, we would choose the model that only uses the *subsidy* as an instrumental variable, as its F-statistic is the highest out of all 3 IV models, even higher than when using both instrumental variables. As the estimates of these two models seem to be the same, we choose the simpler model which only contains the *subsidy* instrumental variable. Note that *distance* appears to be a weak instrument, as its model's F-statistic is low and its estimates are close to the, likely inconsistent, OLS estimates.

Table 6: Regression for log-earnings on schooling.

	<i>Dependent variable:</i>			
	logwage			
	<i>OLS</i>	<i>instrumental variable</i>		
	(1)	(2)	(3)	(4)
schooling	0.101 (0.073)	0.100 (0.626)	0.730*** (0.262)	0.658*** (0.245)
Constant	5.778*** (0.512)	5.789 (4.163)	1.599 (1.753)	2.078 (1.639)
IV variable	-	distance	subsidy	both
F Statistic IV	-	5.6453	41.1994	23.0049
Observations	416	416	416	416
R <sup>2</sup>	0.005	0.005	-0.176	-0.137
Adjusted R <sup>2</sup>	0.002	0.002	-0.179	-0.139
Residual Std. Error (df = 414)	3.467	3.467	3.768	3.705
F Statistic	1.950 (df = 1; 414)			
<i>Note:</i>		*p<0.1; **p<0.05; ***p<0.01		

Table 7: P-values of Hausman test for the IV models

IV	pvalue
IV Distance	0.9979
IV Subsidy	0.0061
IV Both	0.0103

## iii)

The OLS results are similar to the results with the IV regression where we use *distance* as an instrumental variable. The estimates for the other two IV regressions differ strongly from the OLS estimates, but are similar to each other.

If there is no endogeneity, then we would use the OLS model, as it is efficient in this context. Assuming that we have valid instruments, we can test for endogeneity by the Hausman test.

We present the results for the Hausman test in Table 7, clearly we cannot reject the Hausman test for the *distance* instrument variable. However, we can reject the Hausman test on  $\alpha = 0.05$  for the *subsidy* and for both instrumental variables.

These results and the theoretical background suggest, that there is an endogeneity problem and that *distance* is a weak instrument. Hence, we would choose the simplest model which seems to address this endogeneity issue, which is the model with the *subsidy* instrument.