

Introduction to Econometrics

Final exam

The time limit is 90 minutes and the exam is worth a total of 30 points. You are NOT allowed to look up for solutions in books, notes or internet and not allowed to consult the problems with your classmates or anyone knowledgeable. Any violation of academic honesty will be punished to the fullest extent possible.

Date: 15.01.2021

Instructor: Dali Laxton

Multiple choice questions

30 min, 2 points each

1) Consider the following simple regression model $y = \beta_0 + \beta_1 x_1 + u$. The variable z is a poor instrument for x if _____.

- a. there is a high correlation between z and x
- b. there is a low correlation between z and x
- c. there is a high correlation between z and u
- d. there is a low correlation between z and u

2) In the following regression equation, y is a binary variable:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u$$

In this case, the estimated slope coefficient, $\widehat{\beta}_1$ measures _____.

- a. the predicted change in the value of y when x_1 increases by one unit, everything else remaining constant
- b. the predicted change in the value of y when x_1 decreases by one unit, everything else remaining constant
- c. the predicted change in the probability of success when x_1 increases by one unit, everything else remaining constant
- d. the predicted change in the probability of success when x_1 decreases by one unit, everything else remaining constant

3) Which of the following assumptions is required to obtain a first-differenced estimator in a two-period panel data analysis?

- a. The idiosyncratic error at each time period is uncorrelated with the explanatory variables in both time periods.
- b. The explanatory variable does not change over time for any cross-sectional unit.
- c. The explanatory variable changes by the same amount in each time period.
- d. The variance of the error term in the regression model is not constant.

4) Consider the following regression model: $\log(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_{12} + \beta_3 x_3 + u$. This model will suffer from functional form misspecification if _____.

- a. β_0 is omitted from the model
- b. u is heteroskedastic
- c. x_{12} is omitted from the model
- d. x_3 is a binary variable

5) Consider the following regression equation: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u$

In which of the following cases, the dependent variable is binary?

- a. y indicates the gross domestic product of a country
- b. y indicates whether an adult is a college dropout
- c. y indicates household consumption expenditure
- d. y indicates the number of children in a family

Problem 1. (30 min) We want to measure the impact of holding a health insurance (*healthin*) on the medical expenses (*medexp*). The following is the simple model expressing the relationship:

$$\log(\text{medexp}) = \beta_0 + \beta_1 \text{healthin} + \varepsilon$$

- a) (2pt) Why might *healthin* be correlated with ε ?
- b) (2pt) Explain why *healthin* is likely to be related to the *age* and *illnesses* of the insured. Does this mean *age* and *illnesses* are good IV for *healthin*? Why or why not?
- c) (2pt) After controlling for *age* and *illnesses*, you still believe that *healthin* suffers from endogeneity issue. In particular, you believe that the risk-aversion of individuals drives both variables *healthin* and *medexp*. Justify why social security income replacement rate¹ may be a good instrument.
- d) (2pt) How would you proceed with the estimation using the IV? (describe the 2SLS technique in this particular example).
- e) (2pt) Propose an alternative instrument in order to solve the endogeneity issue.

¹ A social security income replacement rate is the percentage of a worker's pre-retirement income that is paid out by a pension program after retirement.

Problem2. (30 min) Suppose you want to assess the impact of a race of an individual on the likelihood of approving a mortgage loan. In the example below the key explanatory variable is *white*, a dummy variable equal to one if the applicant was white. The other applicants in the data set are black and Hispanic. To test for the discrimination in the mortgage loan market, a linear probability model (LPM) can be used:

$$approve = \alpha_0 + \alpha_1 white + u$$

- a) **(2pt)** Suppose you obtain the following output from the regression above. Interpret the coefficient on *white*.

VARIABLES	(1) approve
white	0.201*** (0.0198)
Constant	0.708*** (0.0182)
Observations	1,989
R-squared	0.049
Standard errors in parentheses *** p<0.01, ** p<0.05, * p<0.1	

- b) **(2pt)** Name at least one pro and one con of using an LPM.
 c) **(2pt)** Suppose now that you run probit and logit models as well, interpret the coefficients on *white* for probit and logit models and compare them with the LPM model.

VARIABLES	(LPM) approve	(Probit) approve	(Logit) approve
white	0.201*** (0.0198)	0.784*** (0.0867)	1.409*** (0.151)
Constant	0.708*** (0.0182)	0.547*** (0.0754)	0.885*** (0.125)
Observations	1,989	1,989	1,989
R-squared	0.049		
Standard errors in parentheses *** p<0.01, ** p<0.05, * p<0.1			

- d) **(2pt)** By how much is it more likely for white people to obtain mortgage loan in comparison to minorities according to probit model? How different is this result from LPM result?
 e) **(2pt)** By how much is it more likely for white people to obtain mortgage loan in comparison to minorities according to logit model? Note that the functional form of the logit model is $\Lambda(\cdot) = \frac{\exp(\cdot)}{1+\exp(\cdot)}$.