

# The Exam

**EC 425/525:** Econometrics

Typed, readable submissions due at 5PM on Friday, 14 June 2019

## Instructions

**DUE** Your submission is due at 5PM on Friday, 14 June 2019 on [Canvas](#).

Your submission **must be typed with R code and output/figures beneath your responses**.

E.g., `knitr` and `R Markdown`. Word files are fine.

**IMPORTANT** **You may not work with anyone else**. You may not communicate with anyone about this exam except Jenni and Ed. Any evidence of a violation of this rule will result in a grade of zero.

**IMPORTANT<sub>2</sub>** **All materials are allowed**—notes, previous assignments/solutions, Google, etc..

**IMPORTANT<sub>3</sub>** You may use any functions/packages you would like on this exam—and any statistical package (R, Stata, etc.). I can only guarantee technical support from Jenni and me in R.

**OBJECTIVE** This exam will test your knowledge and help you continue building intuition and skills.

**POINTS** There are 104 points possible on this exam.

**LASTLY** Good luck and thank you for a great quarter.

## Section 1: ATE, TOT, ITT, or LATE

**32 points** (4 points per part)

**Without explanation** Full credit for correct answers; 0 points for incorrect answers.

**With explanation** 50% for correct answer. Up to 50% for incorrect answer with good explanations.

In this section, choose whether the described point estimate is the ATE, TOT, ITT, or LATE. There may be multiple answers; in those cases choose the most general answer. Assume all estimates can be interpreted as causal. Treatment effects can be heterogeneous.

1. Consider the OLS-based estimate of  $\tau$  from the model

$$Y_i = \alpha + \tau D_i + \varepsilon_i$$

Which type of treatment effect does  $\hat{\tau}$  represent?

2. Consider the case of blocking on propensity scores. With  $n$  individuals we calculate

$$\hat{\tau} = \frac{\sum_b n_b \hat{\tau}_b}{n}$$

where  $\hat{\tau}_b$  is the estimated treatment effect in block  $b$  and  $n_b$  is the number of observations in block  $b$ .

Which type of treatment effect does  $\hat{\tau}$  represent?

3. Let  $X_i$  be a continuous variable. When  $X_i$  exceeds the value  $t$ , the probability that individual  $i$  participates in a program increases from 0.2 to 0.4.

We regress our outcome  $Y_i$  on an indicator for  $X_i > t$ . Which type of treatment effect does the coefficient on this indicator represent?

4. Using nearest-neighbor matching on a set of covariates, we match the outcome for each of our treated individuals  $Y_i$  to the outcome for her nearest control neighbor  $Y_j^{(i)}$ . We then calculate the average difference between  $Y_i$  and  $Y_j^{(i)}$ , i.e.,

$$\hat{\tau} = \text{Avg}(Y_i - Y_j^{(i)})$$

Which type of treatment effect does  $\hat{\tau}$  represent?

5. A company approves credit-card applications if and only if the applicant's credit score exceeds 600. To test the effect of credit-card approval on consumption, you estimate the linear model

$$\text{Consumption}_i = \beta_0 + \beta_1 (\text{Score}_i - 600) + \beta_2 \text{Approved}_i + \beta_3 (\text{Score}_i - 600) \times \text{Approved}_i + \varepsilon_i$$

Which type of treatment effect does  $\hat{\beta}_2$  represent?

6. You estimate  $\tau$  via instrumental variables. Which type of treatment effect does  $\hat{\tau}$  represent?

7. You estimate  $\tau$  via instrumental variables **and** there are no always takers. Which type of treatment effect does  $\hat{\tau}$  represent? *Hint:* Your answer should differ from **6**.

8. You regress your outcome  $Y_i$  on an instrument  $Z_i$ . What does the coefficient on  $Z_i$  represent?

## Section 2: The value of a business degree

**42 points** (6 points per part)

Suppose we want to estimate the causal effect of a business degree (**Business<sub>i</sub>**) at UO on earnings 10 years after graduation (**Income<sub>i</sub>**). The university gives you a random sample of undergraduates at UO who applied to be a major in the business school. Assume a homogeneous treatment effect.

For each individual, you observe

- **Income<sub>i</sub>**, the individual's income 10 years after graduation
- **Business<sub>i</sub>**, an indicator for whether the student was accepted by the business school
- **SAT<sub>i</sub>**, the individual's SAT score (from their application to attend UO)
- **GPA<sub>i</sub>**, the individual's GPA at the time of their application to the business school
- **Econ<sub>i</sub>**, an indicator for whether the student passed the introductory economics sequence

**9.** Your first impulse is to regress **Income<sub>i</sub>** on **Business<sub>i</sub>**, **SAT<sub>i</sub>**, **GPA<sub>i</sub>**, and **Econ<sub>i</sub>**. Under which conditions will the coefficient on **Business<sub>i</sub>** be interpretable as the causal effect of a business degree at UO on income?

**10.** You had such a great time coding up the matching estimator on the first problem set in 525, you decide to do it again. Specifically, match admitted individuals to non-admitted individuals who have identical values of **SAT<sub>i</sub>**, **GPA<sub>i</sub>**, and **Econ<sub>i</sub>**. Within each combination of **SAT<sub>i</sub>**, **GPA<sub>i</sub>**, and **Econ<sub>i</sub>**, you calculate the difference between the average income for admitted individual and the average income for rejected individuals.

- a. Under which conditions can we interpret this estimate as a causal effect?
- b. Would you ever prefer the estimate in **9** to this estimate in **10**?

**11.** After several emails, UO offers to provide you with **GPA<sub>i</sub><sup>Final</sup>**, the GPAs of individuals when they graduate from UO. Your contact thinks this variable could be important, as individuals with higher college GPAs tend to make more money.

- a. Are there advantages or disadvantages to replacing **GPA<sub>i</sub>** with **GPA<sub>i</sub><sup>Final</sup>** in **9**?
- b. Are there any advantages or disadvantages to adding this new variable **GPA<sub>i</sub><sup>Final</sup>** as an additional control in the regression described in **9**?

**12.** Suddenly, you remember how cool propensity-score matching seemed. Using a logistic regression, you regress **Business<sub>i</sub>** on **SAT<sub>i</sub>**, **GPA<sub>i</sub>**, and **Econ<sub>i</sub>**. You then add the predictions from this regression ( $\hat{p}_i$ ) to the regression described in **9**.

Under which conditions will this estimate be interpretable as causal?

**13.** During a conversation with an administrator from the business school, you learn that the business only admits students if two conditions are true: **(1)** their GPA is above 3.0, and **(2)** the student passed the introductory economics sequence.

Describe an estimation strategy in which you could credibly estimate the causal effect of admittance into the business school using the variables described at the start of this section. Include any assumptions your strategy makes.

## Section 3: A simulation (of course!)

**30 points** (5 points per part)

Back to instrumental variables. We're going to explore some violations/bastardizations of IV/2SLS.

### The data-generating process

Set up a data-generating process

$$Y_i = (\beta_0 = 3) + (\beta_1 = 2) X_i + u_i \quad (1)$$

where  $X_i$  is correlated with  $u_i$ . In your data-generating process, include a valid instrument  $Z_i$  for  $X_i$ . Set up the variance-covariance matrix  $\Sigma$  among the variables  $X_i$ ,  $u_i$ , and  $Z_i$

$$\Sigma = \text{Cov}(X_i, u_i, Z_i) = \begin{bmatrix} 1 & 0.5 & 0.7 \\ 0.5 & 1 & 0 \\ 0.7 & 0 & 1 \end{bmatrix}$$

Each variable is mean zero and comes from a multivariate Normal distribution. For each iteration (run at least 1,000 iterations), sample  $N = 50$  individuals from this defined DGP.

Finally, generate the variable  $W_i = X_i + \varepsilon_i$  where  $\varepsilon_i \sim N(0, 1)$ .

### The models

In each iteration, estimate the effect of  $X_i$  on  $Y_i$  using each of the following five models:

- A.** The OLS estimate for  $\beta_1$  in equation (1)
- B.** The 2SLS estimate of  $\beta_1$  in equation (1), in which  $Z_i$  instruments for  $X_i$
- C.** The 2SLS estimate of  $\beta_1$  in equation (1), in which the indicator  $\mathbb{I}\{Z_i > 0\}$  instruments for  $X_i$
- D.** The OLS estimate of  $\beta_1$  in (1) when we add  $Z_i$  as a control.
- E.** The 2SLS estimate of  $\beta_1$  in (1) using  $W_i$  as a control (in both stages).

### The actual questions

- 14.** Which elements of  $\Sigma$  tell us that  $Z_i$  is a valid instrument?
- 15.** Which models appear to provide unbiased estimates? Provide graphical and numeric summaries of each model (with labeled axes and models clearly distinguished).
- 16.** Which models do we expect to consistently estimate the causal effect of  $X_i$  on  $Y_i$ ? Explain why.
- 17.** How does the model described in **D.** differ from the typical IV/2SLS setup? Explain how this difference could lead to biased/inconsistent estimates for  $\beta_1$ .
- 18.** What is wrong with the model described in **E.**? Explain.
- 19.** In this setting, does the indicator variable  $\mathbb{I}\{Z_i > 0\}$  provide a valid instrument? How does it affect efficiency, relative to the 2SLS model described **B.**? Explain.

**The end!**