

Week 2: OVB

JJ Chen

January 23, 2015

Anouncement: Today's Tasks

- ▶ OVB:
 - ▶ Regression Anatomy and OVB formula
 - ▶ Typical Usage of OVB formula
 - ▶ Bad controls
- ▶ Examples:
 - ▶ *Estimating the payoff to attending a more selective college: An application of selection on observables and unobservables* (Dale and Krueger (2002))
 - ▶ *Selection on Observed and Unobserved Variables: Assessing the Effectiveness of Catholic Schools* (Altonji, Elder, and Taber (2005))

Week 2: OVB

Regression Anatomy (the FWL Theorem)

- ▶ Consider a bivariate population regression: $Y_i = \alpha + \beta D_i + e_i$
 - ▶ take covariance in either side,

$$\text{Cov}(Y_i, D_i) = \text{Cov}(\alpha + \beta D_i + e_i, D_i) \implies$$
 - ▶ $\beta = \frac{\text{Cov}(Y_i, D_i)}{\text{Var}(D_i)}$ (how about α ?)
- ▶ Consider a multivariate population regression:

$$Y_i = \alpha + \beta D_i + X_i' \gamma + e_i$$
 - ▶ regression anatomy; how to get β ?
 - ▶ Step 1 (optional): run a regression of Y_i on X_i and get the residual \tilde{Y}_i
 - ▶ Step 2: run a regression of D_i on X_i and get the residual \tilde{D}_i
 - ▶ Step 3: run a bivariate regression of \tilde{Y}_i on \tilde{D}_i
- ▶ Thus $\beta = \frac{\text{Cov}(\tilde{Y}_i, \tilde{D}_i)}{\text{Var}(\tilde{D}_i)}$, or $\beta = \frac{\text{Cov}(Y_i, \tilde{D}_i)}{\text{Var}(\tilde{D}_i)}$ (Why?)

OVB Formula

- ▶ Consider a short and a long regression:
 - ▶ $Y_i = \alpha^s + \beta^s D_i + e_i^s$
 - ▶ $Y_i = \alpha^l + \beta^l D_i + \gamma^l X_i + e_i^l$
- ▶ Taking covariance for the short regression and substituting the long regression gives a simple OVB formula

$$\beta^s = \frac{\text{Cov}(Y_i, D_i)}{\text{Var}(D_i)} \implies$$

$$\beta^s = \frac{\text{Cov}(\alpha^l + \beta^l D_i + \gamma^l X_i + e_i^l, D_i)}{\text{Var}(D_i)} \implies$$

$$\beta^s = \beta^l + \gamma^l \frac{\text{Cov}(X_i, D_i)}{\text{Var}(D_i)} = \beta^l + \gamma^l \pi_{XD}$$

- ▶ where π_{XD} is from the auxiliary regression $X_i = \phi + \pi_{XD} D_i + u_i$
- ▶ short = long + OVB = long + the effect of omitted \times the regression of omitted on include (what if more than one OV?)

Typical Usage of the OVB Formula

- ▶ Often times, omitted variables are variables we wish we had
- ▶ OVB formula is thus used to predict the likely direction of bias
 - ▶ what's the sign of ability bias in a wage-schooling regression?
 - ▶ ability is positively correlated with earning
 - ▶ theory predicts that correlation between ability and schooling can be ambiguous (why?)
- ▶ More exercises:
 - ▶ $\text{sex}_{ir} = \alpha + \rho \text{HIV}_r + \mathbf{X}'_{ir} \boldsymbol{\gamma} + u_i$
 - ▶ $\text{obesity}_{ir} = \alpha + \rho \text{presidential}_s \text{prawl}_r + \mathbf{X}'_{ir} \boldsymbol{\gamma} + u_i$
- ▶ Some time we want to estimate the extent of OVB
 - ▶ how large is the OVB if we want to explain away the entire estimate of your β^s
 - ▶ catholic school example

Bad Controls: What Should We Control for

- ▶ Recall from CIA, a good control X_i allows us to say:
 - ▶ conditional on X_i , the assignment of D_i is as good as randomly assigned
 - ▶ $E(Y_0 | D_i, X_i) = E(Y_0 | X_i)$
 - ▶ $E(\eta_i | D_i, X_i) = E(\eta_i | X_i)$
- ▶ By the logic of OVB, it seems that we should control for all variables X_i that are correlated with both the outcome variable Y_i and the treatment D_i
 - ▶ but we also said that X_i should be pre-treatment variables
 - ▶ what happened if we controls for post-treatment variables (or intermediate outcomes)?

Bad Controls: Examples

- ▶ Example: suppose we want to know the impact of parents' occupations on kids' test scores: $Y_i = \alpha + \beta \text{occu}_i + e_i$
 - ▶ should we control for family incomes (potentially affected by occupations)?
 - ▶ is it a good control in terms of giving $E(\eta_i \mid \text{occu}, \text{income}) = E(\eta_i \mid \text{income})$?
 - ▶ does it make sense to compare two people with the same income (say 200K) but different occupations? (give an example)
- ▶ Example from MHE: should we control for occupations in a wage-schooling regression?
 - ▶ occupation is correlated with both the treatment (years of schooling) and outcome (earning)

Table 1: Bad controls and selection bias

- ▶ ATE: 0.5
- ▶ Controlling for occupations: 0

The Question and Selection Problem

- ▶ Let's discuss Dale and Kruger's paper and review basic concepts in this week
- ▶ Question: "Does the 'quality' of the college that students attend influence their subsequent earnings?"
- ▶ Selection Problem:
 - ▶ Can naive comparisons between students from selective colleges and non-selective colleges answer the question?
 - ▶ What is the possible selection bias?
 - ▶ If we control for family income, gender, race, SAT scores, and high school ranking, will the selection bias disappear?
 - ▶ If not, what will drive the bias term and how will you predict the sign of the bias?

Solutions to the Selection Problem

- ▶ Choices made by school admission committees and students in the application, screen, and match process might reveal useful information
 - ▶ when applying for schools, students reveal their potential ability by the choice of schools they apply to
 - ▶ admission committees have the chance to read students' essays, letters of recommendation, etc
 - ▶ given application outcomes, the decision of attending a particular school again reveals unobserved characteristics
- ▶ Extend “selection on the observables” to “selection on the observables and unobservables”
 - ▶ “match-applicant model”
 - ▶ “self-revelation model”

References I

Altonji, Joseph G, Todd E Elder, and Christopher R Taber. 2005. "Selection on Observed and Unobserved Variables: Assessing the Effectiveness of Catholic Schools." *Journal of Political Economy* 113 (1). ERIC: 151.

Dale, Stacy Berg, and Alan B Krueger. 2002. "Estimating the Payoff to Attending a More Selective College: An Application of Selection on Observables and Unobservables*." *The Quarterly Journal of Economics* 117 (4). MIT Press: 1491–1527.