

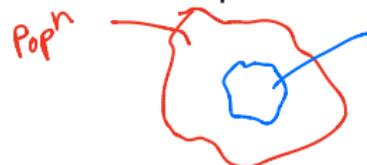
Tutorial 12: Sample Selection

Hammad Shaikh

April 8, 2021

Sample Selection Introduction

- Empirical analysis based on samples and not the population



↳ Empirics use sample to infer about population

- Sample used for analysis may not represent the population



↳ Empirics may not be informative
about males

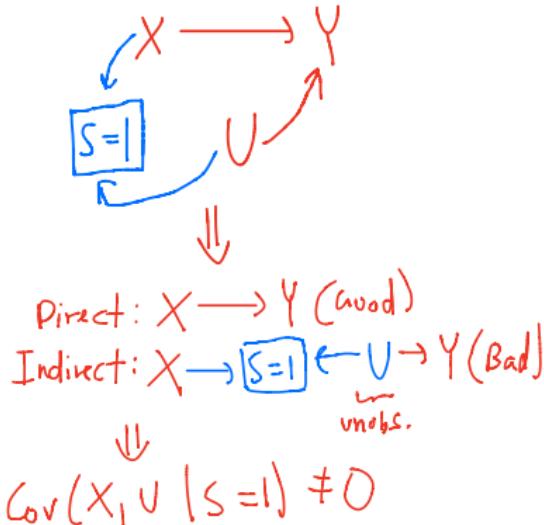
- Participants may select themselves to the sample of interest



↳ participant is choosing to be
employed

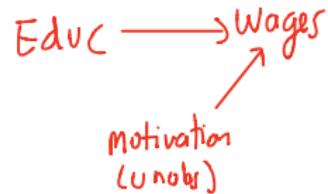
} Focus of
today

Formalization of the Selection Problem

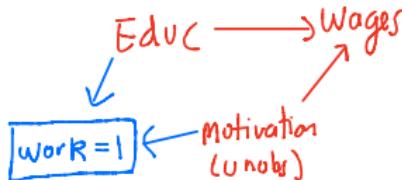


- Population: $Y_i = \alpha_0 + \alpha_1 X_i + U_i, E(U_i | X_i) = 0$
 - $\hookrightarrow \{(X_i, Y_i)\}_{i=1}^N$ where $N = \text{population size}$
 - $\hookrightarrow \text{observe everyone}$
- Sample selection: $s_i = I(Z_i \gamma_1 + X_i \gamma_2 + V_i \geq 0)$
 - $\hookrightarrow S_i = \begin{cases} 1, & (X_i, Y_i) \text{ obs.} \\ 0, & \text{do not obs. } Y_i \end{cases}$ Exog. Educ. Unobs.
 - \star Selection regt. charac. who is in sample and who is not
- Conditional expectation for selected sample: $E(Y_i | X_i, Z_i, s_i = 1) = \beta_0 + X_i \beta_1 + E(u_i | X_i, Z_i, s_i = 1)$
 - $\underbrace{\text{in sample}}$
 - $\underbrace{g \times \frac{\varphi(z_i \gamma_1 + X_i \gamma_2)}{\Phi(z_i \gamma_1 + X_i \gamma_2)} \}_{\text{Inverse Mills Ratio}}$
 - $\Rightarrow X_i \text{ is correlated with } E[u_i | X_i, t_i, s_i = 1]$
- Selection bias: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ if $\underbrace{\text{Cov}(X_i, \epsilon_i | s_i = 1)}_{\text{contain inverse mills ratio (omitted)}} \neq 0$
 - \hookrightarrow omitted variable bias

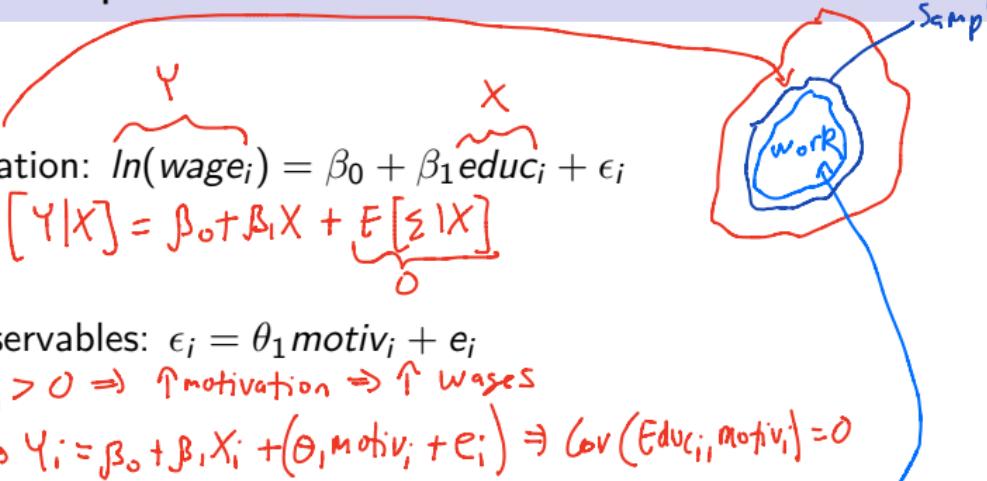
Example of Sample Selection Bias



* If obs. popn, OLS unbiased & consistent



* In sample we get a omitted variable bias since Educ & motivation correlated through sample selection.

- 
- Population: $\ln(wage_i) = \beta_0 + \beta_1 \text{educ}_i + \epsilon_i$
 $\hookrightarrow E[Y|X] = \beta_0 + \beta_1 X + E[\epsilon|X]$
 - Unobservables: $\epsilon_i = \theta_1 \text{motiv}_i + \eta_i$
 $\hookrightarrow \theta_1 > 0 \Rightarrow \uparrow \text{motivation} \Rightarrow \uparrow \text{wages}$
 $\hookrightarrow Y_i = \beta_0 + \beta_1 X_i + (\theta_1 \text{motiv}_i + \eta_i) \Rightarrow \text{Cov}(\text{Educ}_i, \text{motiv}_i) \neq 0$
 - Sample selection: $\text{work}_i = I(\gamma \text{educ}_i + (\theta_2 \text{motiv}_i + \eta_i) \geq 0)$
 $\hookrightarrow \gamma > 0, \theta_2 > 0$
 $\rightarrow \text{Cov}(\Sigma_i, V_i) \neq 0 \Rightarrow \text{Selection problem}$
 - Selection bias: $\text{Cov}(\text{educ}_i, \text{motiv}_i | \text{work}_i = 1) < 0$
 $\hookrightarrow \text{Suppose } \gamma \text{ is given } \downarrow \text{educ}_i \text{ & } \text{work}_i = 1$
 $\Rightarrow \uparrow \text{motiv}_i \Rightarrow E[\underbrace{\text{motiv} \times \text{educ}}_{\text{unobs.}} | \text{work}_i = 1] \neq 0$

↳ Selection bias for $\log(wage_i) = \beta_0 + \beta_1 \text{Educ}_i + \beta_2 \text{motiv}_i$ only for $\text{work}_i = 1$

Solutions to Selection Problem

motivation

control
variable

→ ①

Include control variable W_i so that $\text{Cov}(X_i, \epsilon_i, W_i | s_i = 1) = 0$

↳ $\log(wage_i) = \beta_0 + \beta_1 \text{educ}_i + \beta_2 \underbrace{\text{motiv}_i}_{\text{observed}} + \epsilon_i$

↳ No endog. issue since $\text{Cov}(\text{educ}_i, \epsilon_i | \text{motiv}_i) = 0$

control
function

→ ②

Control for the estimated Inverse Mills Ratio (Heckman)

↳ (i) Estimate the selection eqn (eg. probit)

$$\hat{\delta}_i = I(z_i \hat{r}_1 + x_i \hat{r}_2)$$

↳ (ii) Control for estimated Inverse Mills Ratio

$$\frac{\phi(z_i \hat{r}_1 + x_i \hat{X}_2)}{\Phi(z_i \hat{r}_1 + x_i \hat{X}_2)} = \hat{\lambda}_i(z_i, x_i)$$

$$Y_i = \beta_0 + \beta_1 X_i + p \cdot \hat{\lambda}_i + \epsilon_i$$

Implement
in Stata
together

control for inverse mills ratio