# Recitation 10

Seung-hun Lee

Columbia University

## Note on the final project

- Approach the TA's for advice: regression methods, data, or ideas

- Also, refer to my recitation notes for Diff-in-Diffs and Regression discontinuity methods - I explain how you can use it in a natural experiment context

- Data: Household surveys are the safe place to start but country level data could work depending on your idea

- But first, **be clear about what kind of question you want to explore**: are you looking for a correlation/causal relation? Are you trying to predict something?

- Who to look for (topic-wise) - all three of us cover basic grounds in most topics, but for specialized stuff...
  - Me: Development/Public/Labor econ (applied microeconomics in general). If you are trying to study health-care markets in the US, I am the wrong person to ask

# Experiments in Economics

Motivation

- You categorize some individuals under treatment group and controlled group and compare the differences across groups and times.
- You can do this in lab setting (randomized control trials) or find an exogenous event (natural experiments)
- They provide a conceptual benchmark for assessing observational studies and can solve many validity threats in regular regressions.
- Do note that they have a validity threat of their own

## Potential Outcome Framework

Setup

- Start by assuming that the treatment effect is identical for everyone (**constant treatment effect** assumption)
- Let $Y_{i,t}$ be the **observed** outcome variable for individual $i$ at time $t$
- $X_{it}$ be the treatment variable: It is 1 if individual $i$ is treated and 0 if otherwise.
- Let $Y_{it}(0)$ denote a **potential** outcome if subject $i$ is not treated at time $t$ and $Y_{it}(1)$ be the same if $i$ is treated.
- Then $Y_{it}$ can be split into

$$Y_{it} = Y_{it}(1)X_{it} + Y_{it}(0)(1 - X_{it})$$
$$= Y_{it}(0) + (Y_{it}(1) - Y_{it}(0))X_{it}$$

# Potential Outcome Framework

Two takeaways

- **Fundamental problem of missing data**: We can only observe at most one of $Y_{it}(1)$ and $Y_{it}(0)$, since individual $i$ cannot be treated and untreated simultaneously at time $t$
  - For us to make statements about the treatment effect, we need to be sure about what the missing outcome looks like.
  - Perfect randomization, randomization conditional on observables, instrumental variables on treatment variable $X_{it}$...
- **Average treatment effect**: Average treatment effect is defined as

$$ATE = E[Y_{it}(1) - Y_{it}(0)] = \frac{1}{N} \sum_{i=1}^{N} (Y_{it}(1) - Y_{it}(0))$$

which is obtained from the coefficient of the $X_{it}$ variable of the regression

$$Y_{it} = \beta_0 + \beta_1 X_{it} + u_{it}$$

# Average Treatment Effect

Regressional form

- We will assume that the experiment is **perfectly randomized**
  - The treated individuals and controlled individuals are identical except for treatment status, or that $E[u_{it}|X_{it}] = 0$ for all possible $X_{it}$ values.
- Derive the expected value of $Y_{it}$ separately for the treated and controlled individuals. For the controlled ($X_{it} = 0$, so that $Y_{it} = Y_{it}(0)$):

$$E[Y_{it}|X_{it} = 0] = E[Y_{it}(0)] = \beta_0$$

and for the treated, ($X_{it} = 1$, so that $Y_{it} = Y_{it}(1)$)

$$E[Y_{it}|X_{it} = 1] = E[Y_{it}(1)] = \beta_0 + \beta_1$$

# Average Treatment Effect

Regressional form

- Then, the average treatment effect can be characterized as

$$ATE = E[Y_{it}(1) - Y_{it}(0)] = (\beta_0 + \beta_1) - \beta_0 = \beta_1$$

where subtraction among $E[Y_{it}|X_{it} = 1]$ and $E[Y_{it}|X_{it} = 0]$ is possible since we are assuming perfect randomization.

- Under perfect randomization, identifying the average treatment effect is equivalent to obtaining the $\beta_1$ coefficient through an OLS process.

- However this is possible in rare circumstances...

# Average Treatment Effect: No constant treatment effects

Approach

- Define $\beta_{1i} = Y_{it}(1) - Y_{it}(0)$ and let $\beta_1 = E[\beta_{1i}]$.
- The potential outcome framework can be formally written as

$$
\begin{aligned}
Y_i &= Y_i(1)X_i + Y_i(0)(1 - X_i) \\
&= Y_i(0) + (Y_i(1) - Y_i(0))X_i \\
&= E[Y_i(0)] + (Y_i(1) - Y_i(0))X_i + Y_i(0) - E[Y_i(0)] \\
&= \beta_0 + \beta_{1i}X_i + u_i
\end{aligned}
$$

- Furthermore,

$$
Y_i = \beta_0 + \beta_1 X_i + \underbrace{(\beta_{1i} - \beta_1)X_i + u_i}_{=v_i}
$$

# Average Treatment Effect: No constant treatment effects

Approach

- As long as we have perfect randomization, even $E[v_i|X_i] = 0$ will hold.

$$
\begin{aligned}
E[v_i|X_i] &= E[(\beta_{1i} - \beta_1)X_i + u_i|X_i] \\
&= E[(\beta_{1i} - \beta_1)X_i|X_i] + E[u_i|X_i] \\
&= E[(\beta_{1i} - \beta_1)|X_i]X_i + E[u_i|X_i] \\
&= 0
\end{aligned}
$$

- With perfect randomization, OLS gives an unbiased estimate of the average treatment effect, whether we assume constant treatment effects or not.

# Validity Threats

Validity threats

- We cannot automatically subtract $E[Y_{it}|X_{it} = 1]$ and $E[Y_{it}|X_{it} = 0]$ under **imperfect randomization**.

- If the **attrition rate is nonrandom** - if it differs by a treatment status - we end up with a biased estimate of the treatment effect.

- In terms of experimental procedure, **failure to comply to experiment protocol** and other **experimental effects** that rises through peculiar behaviors of the experimental and the experiment subject could also serve as a validity threat.

- Constant treatment effect assumption that we have imposed on ourselves may not be accurate.

- There are external validity threats when the sample is not representative, is not replicable in other settings, and the results may have general equilibrium effects.

# Validity Threats

Addressing validity threats

- If the treated and the controlled differ in some observed characteristics, we can simply include control variables $Z_{it}$
  - This gives us ATE conditional on observables and identify treatment effects for people with different treatment effects
- Instrumental variable approach: If there exists a $Z_{it}$ variable that influences the treatment status $X_{it}$ and uncorrelated with $u_{it}$, we can use $Z_{it}$ to instrument $X_{it}$ and run 2SLS regression
  - The predicted value of $X_{it}$: The probability of being treated.
  - 2SLS estimates the causal effect for those whose value of $X_{it}$ is influenced by $Z_{it}$, putting more weight on those more likely to be treated.
  - This effectively identifies the treatment effect 'localized' for those more likely to be treated
  - We call this **localized average treatment effect** (LATE).

# IV in treatment effects

Regression

- Let $Z_{it}$ be an instrument for $X_{it}$ and apply 2SLS method in this manner

$$X_{it} = \pi_0 + \pi_{1i}Z_{it} + e_{it} \text{ (First stage)}$$
$$Y_{it} = \beta_0 + \beta_{1i}X_{it} + u_{it} \text{ (Equation of interst)}$$

- Obtain $\hat{X}_{it} = \hat{\pi}_0 + \hat{\pi}_{1i}X_{it}$ - the predicted probability of being treated.
- Put $\hat{X}_{it}$ in place of $X_{it}$ in the equation of interest.
- If $\beta_{1i}, \pi_{1i}$ are independent of $(u_{it}, v_{it}, Z_{it})$, $E(\pi_{1i}) \neq 0$, then

$$\hat{\beta}_1 \xrightarrow{p} \frac{E(\beta_{1i}\pi_{1i})}{E(\pi_{1i})} \text{ (Refer to appendix 13.2)}$$

- Ultimately, $\hat{\beta}_1$ can be written as

$$\hat{\beta}_1 \xrightarrow{p} E[\beta_{1i}] + \frac{cov(\beta_{1i}, \pi_{1i})}{E(\pi_{1i})} = \beta_1 + \frac{cov(\beta_{1i}, \pi_{1i})}{E(\pi_{1i})}$$

- Treatment effect is large for individuals for whom the effect of the instrument is large.