

# Introduction to Econometrics

## *Lecture 1 : Causal Inference in Social Science*

**Zhaopeng Qu**

**Business School, Nanjing University**

Sep. 13th, 2018



# Outlines

- 1 Causal Inference in Social Science
- 2 Experimental Design as an Benchmark
- 3 Program Evaluation Econometrics

# Outlines

- 1 Causal Inference in Social Science
- 2 Experimental Design as an Benchmark
- 3 Program Evaluation Econometrics

# Outlines

- 1 Causal Inference in Social Science
- 2 Experimental Design as an Benchmark
- 3 Program Evaluation Econometrics

# Causal Inference in Social Science

# The Purposes of Empirical Work

- “The objective of science is the discovery of the relations” — Lord Kelvin
- In most cases, we often want to explore the relationship between two variables in one paper.
  - eg. education and wage
- Then, in simplicity, there are two relationships between two variables.
  - Correlation(相关) V.S. Causality (因果)

## A Classical Example: Hemline Index (裙边指数)

- George Taylor, an economist in the United States, made up the phrase it in the 1920s. The phrase is derived from the idea that hemlines on skirts are shorter or longer depending on the economy.
  - Before 1930s, fashion women favored middle skirts most.
  - In 1929, long skirts became popular. While the *Dow Jones Industrial Index(DJII)* plunged from about 400 to 200 and to 40 two years later.
  - In 1960s, DJII rushed to 1000. At the same time, short skirts showed up.
  - In 1970s, DJII fell to 590 and women began to wear long skirts again.
  - In 1990s, mini skirt debuted, DJII rushed to 10000.
  - In 2000s, bikini became a nice choice for girls, DJII was high up to 13000.
  - So what is about now? Long skirt is resorting?

# Hemline Index:1920s-2010s





# The Core of Empirical Studies: Causality v.s. Forecasting

- Some Big Data researchers think causality is not important any more in our times..
  - “Look at correlations. Look at the ‘what’ rather than the ‘why’, because that is often good enough.” -Viktor Mayer-Schonberger(2013)
- Most empirical economists think that correlation only tell us the superficial, even false relationship while causal relationship can provide solid evidence to make interference to the real relationship.
  - Today, empirical economists care more about the causal relationship of their interests than ever before.
  - “the most interesting and challenging research in social science is about cause and effect” ——Angrist and Lavy(2008)

# The Core of Empirical Studies: Causality v.s. Forecasting

- Machine learning is a set of data-driven algorithms that use data to predict or classify some variable  $Y$  as a function of other variables  $X$ .
  - There are many machine learning algorithm. The best methods vary with the particular data application
- Machine learning is all about prediction.
  - Having a good prediction does work sometimes but does NOT mean understanding causality.

# The Core of Empirical Studies: Causality v.s. Forecasting

- Even though forecasting need not involve causal relationships, economic theory suggests patterns and relationships that might be useful for forecasting.
  - Econometric analysis((times series) allows us to quantify historical relationships suggested by economic theory, to check whether those relationships have been stable over time, to make quantitative forecasts about the future, and to assess the accuracy of those forecasts.
- the biggest difference between **machine learning** and **econometrics(or causal inference)**.

# The Central Question of Causality(I)

- A simple example: Do hospitals make people healthier? (Q: **Dependent variable and Independent variable?**)
- A naive solution: compare the health status of those who have been to the hospital to the health of those who have not.
- Two key questions are documented by the questionnaires (问卷) from *The National Health Interview Survey(NHIS)*
  - ① “During the past 12 months, was the respondent a patient in a hospital overnight?”
  - ② “Would you say your health in general is excellent, very good, good ,fair and poor” and scale it from the number “1” to “5” respectively.

# The Central Question of Causality(II)

## Hospital v.s. No Hospital

<i>Group</i>	Sample Size	Mean Health Status	Std.Dev
<i>Hospital</i>	7774	2.79	0.014
<i>No Hospital</i>	90049	2.07	0.003

- In favor of the non-hospitalized, WHY?
  - Hospitals not only cure but also hurt people.
    - ① hospitals are full of other sick people who might infect us
    - ② dangerous machines and chemicals that might hurt us.
  - More important : people having worse health tends to visit hospitals.
- This simple case exhibits that it is not easy to answer an causal question, so let us formalize an model to show where the problem is.

# The Central Question of Causality(III)

- So A right way to answer a causal questions is construct a counterfactual world, thus “What If ....then” , Such as
- An example: How much wage premium you can get from college attendance(上大学使工资增加多少 ?)
  - For any worker, we want to compare
    - Wage if he have a college degree
    - Wage if he had not a college degree
  - Then make a difference. This is the right answer to our question.

# Difficulty in Identification

- Others are the same as
  - Military service
  - Migration
  - Road building
  - Job training
  - Party membership
  - Public policies
  - Others...
- Difficulty: only one state can be observed

# Formalization: Rubin Causal Model

- Treatment :  $D_i$  is a **dummy** that indicate whether individual  $i$  receive treatment or not

$$D_i = \begin{cases} 1 & \text{if individual } i \text{ received the treatment} \\ 0 & \text{otherwise} \end{cases}$$

- Examples:
  - Go to college or not
  - Have health insurance or not
  - Join a training program or not
  - Make an online-advertisement or not
  - ....



# Formalization: Treatment

- Treatment :  $D_i$  can be a **multiple valued**(countinuous) variable

$$D_i = s$$

- Examples:
  - Schooling years
  - Number of Children
  - Number of advertisements
  - Money Supply
- For simplicity, we assume treatment variable  $D_i$  is just a **dummy**.

# Formalization: Potential Outcomes

- A potential outcome is the outcome that would be realized if the individual received a specific value of the treatment.
  - Annual earnings if attending to college
  - Annual earnings if not attending to college
- For each individual, we have two potential outcomes,  $Y_{1i}$  and  $Y_{0i}$ , one for each value of the treatment
  - $Y_{1i}$  : Potential outcome for an individual  $i$  with treatment.
  - $Y_{0i}$  : Potential outcome for an individual  $i$  with treatment.

$$\text{Potential Outcomes} = \begin{cases} Y_{1i} & \text{if } D_i = 1 \\ Y_{0i} & \text{if } D_i = 0 \end{cases}$$

# Stable Unit Treatment Value Assumption (SUTVA)

- Observed outcomes are realized as

$$Y_i = Y_{1i}D_i + Y_{0i}(1 - D_i)$$

- Implies that potential outcomes for an individual  $i$  are unaffected by the treatment status of other individual  $j$
- Individual  $j$ 's potential outcomes are only affected by his/her own treatment.
- Rules out possible treatment effect from other individuals (**spillover effect/externality**)
  - Contagion
  - Displacement

## Causal effect for an Individual

- To know the difference between  $Y_{1i}$  and  $Y_{0i}$ , which can be said to be the **causal effect** of going to college for individual  $i$ . (Do you agree with it?)

### Definition

**Causal inference** is the process of estimating a comparison of counterfactuals under different treatment conditions on the same set of units. It also call **Individual Treatment Effect(ICE)**

$$\delta_i = Y_{1i} - Y_{0i}$$

# Formalization: Estimate ICE

- Due to unobserved counterfactual outcome, we need to make strong assumptions to estimate ICE.
  - Rule out that the ICE differs across individuals ( “heterogeneity effect” )
- Knowing individual effect is not our final goal. As a social scientist, we would like more to know the **Average** effect as a **social pattern**.
- So it make us focus on the average wage for a group of people.
  - How can we get the average wage premium for college attendance?

# Conditional Expectation:

- **Expectation:** We usually use  $E[Y_i]$  (the expectation of a variable  $Y_i$ ) to denote population average of  $Y_i$ 
  - Suppose we have a population with  $N$  individuals

$$E[Y_i] = \frac{1}{N} \sum_{i=1}^N Y_i$$

- **Conditional Expectation:**

- The average wage for those who attend college:  $E[Y_i | D_i = 1]$
- The average wage for those who did not attend college:  $E[Y_i | D_i = 0]$

# Average Causal Effects

## Average Treatment Effect (ATE)

$$\alpha_{ATE} = E[\delta_i] = E[Y_{1i} - Y_{0i}]$$

- It is average of ICEs over **the population**.

## Average treatment effect on the treated(ATT)

$$\alpha_{ATT} = E[\delta_i | D_i = 1] = E[Y_{1i} - Y_{0i} | D_i = 1]$$

- Average of ICEs over the **treated population**

# Fundamental Problem of Causal Inference

- We can never directly observe causal effects (ICE, ATE or ATT)
- Because we can never observe both potential outcomes ( $Y_{0i}$ ,  $Y_{1i}$ ) for any individual.
- We need to compare **potential outcomes**, but we only have **observed outcomes**
- So by this view, causal inference is a **missing data** problem.



# Fundamental Problem of Causal Inference

- Imagine a population with 4 people

i	$Y_{i1}$	$Y_{0i}$	$Y_i$	$D_i$	$Y_{i1} - Y_{0i}$
Tom	3	?	3	1	?
Jerry	2	?	2	1	?
Scarlett	?	1	1	0	?
Nicole	?	1	1	0	?

- What is Individual causal effect (ICE) of attending college for Tom? for Nicole?

# Individual Causal Effect

- Suppose we can observe counterfactual outcomes

i	$Y_{i1}$	$Y_{0i}$	$Y_i$	$D_i$	$Y_{i1} - Y_{0i}$
Tom	3	2	3	1	1
Jerry	2	1	2	1	1
Scarlett	1	1	1	0	0
Nicole	1	1	1	0	0

- The ICE for Tom

$$\delta_{Tom} = 3 - 2 = 1$$

- The ICE for Nicole

$$\delta_{Nicole} = 1 - 1 = 0$$

# Average Treatment Effect(ATE)

- Missing data problem also arises when we estimate ATE

i	$Y_{1i}$	$Y_{0i}$	$Y_i$	$D_i$	$Y_{1i} - Y_{0i}$
Tom	3	?	3	1	?
Jerry	2	?	2	1	?
Scarlett	?	1	1	0	?
Nicole	?	1	1	0	?
$E[Y_{1i}]$	?				
$E[Y_{0i}]$		?			
$E[Y_{1i} - Y_{0i}]$					?

- What is the effect of attending college on average wage of population(ATE)

$$\alpha_{ATE} = E[\delta_i] = E[Y_{1i} - Y_{0i}]$$

# Average Treatment Effect(ATE)

- Missing data problem also arises when we estimate ATE

i	$Y_{1i}$	$Y_{0i}$	$Y_i$	$D_i$	$Y_{1i} - Y_{0i}$
Tom	3	2	3	1	1
Jerry	2	1	2	1	1
Scarlett	1	1	1	0	0
Nicole	1	1	1	0	0
$E[Y_{1i}]$	$\frac{3+2+1+1}{4} = 1.75$				
$E[Y_{0i}]$	$\frac{2+1+1+1}{4} = 1.25$				
$E[Y_{1i} - Y_{0i}]$	0.5				

- What is the effect of attending college on average wage of *the population*(ATE)

$$\alpha_{ATE} = E[\delta_i] = E[Y_{1i} - Y_{0i}] = \frac{1 + 1 + 0 + 0}{4} = 0.5$$

# Average Treatment Effect on the Treated(ATT)

- Missing data problem arises when we estimate ATT

i	$Y_{1i}$	$Y_{0i}$	$Y_i$	$D_i$	$Y_{1i} - Y_{0i}$
Tom	3	?	3	1	?
Jerry	2	?	2	1	?
Scarlett	?	1	1	0	?
Nicole	?	1	1	0	?
$E[Y_{1i} D_i = 1]$	?				
$E[Y_{0i} D_i = 1]$		?			
$E[Y_{1i} - Y_{0i} D_i = 1]$					?

- What is the effect of attending college on average wage for *those who attend college*(ATT)

$$\alpha_{ATE} = E[\delta_i] = E[Y_{1i} - Y_{0i}|D_i = 1]$$

# Average Treatment Effect on the Treated(ATT)

- Missing data problem also arises when we estimate ATE

i	$Y_{1i}$	$Y_{0i}$	$Y_i$	$D_i$	$Y_{1i} - Y_{0i}$
Tom	3	2	3	1	1
Jerry	2	1	2	1	1
Scarlett	1	1	1	0	0
Nicole	1	1	1	0	0
$E[Y_{1i} D_i = 1]$	$\frac{3+2}{2} = 2.5$				
$E[Y_{0i} D_i = 1]$	$\frac{2+1}{2} = 1.5$				
$E[Y_{1i} - Y_{0i} D_i = 1]$	1				

- The effect of attending college on average wage for *those who attend college*(ATT)

$$\alpha_{ATE} = E[Y_{1i} - Y_{0i}|D_i = 1] = \frac{1 + 1}{2} = 1$$

# Observed Association and Selection Bias

- Causality is defined by **potential outcomes**, not by **realized (observed) outcomes**.
- In fact, we can not observe all potential outcomes .Therefore, we can not estimate the above causal effects without further assumptions.
- By using observed data, we can only establish **association (correlation)**, which is the observed difference in average outcome between those getting treatment and those not getting treatment.

$$\alpha_{corr} = E[Y_{1i}|D_i = 1] - E[Y_{0i}|D_i = 0]$$

# College vs Non-College Wage Differentials:

- Comparing the average wage in labor market who went to college and did not go.

## College vs Non-College Wage Differentials:

$$\begin{aligned}
 &= E[Y_{1i}|D_i = 1] - E[Y_{0i}|D_i = 0] \\
 &= \{E[Y_{1i}|D_i = 1] - E[Y_{0i}|D_i = 1]\} + \{E[Y_{0i}|D_i = 1] - E[Y_{0i}|D_i = 0]\}
 \end{aligned}$$

- Question 1: Which one defines the causal effect of college attendance?



# Formalization: Rubin Causal Model

- **Selection Bias(SB)** implies the potential outcomes of treatment and control groups are different even if both groups receive the same treatment

$$E[Y_{0i}|D_i = 1] - E[Y_{0i}|D_i = 0]$$

- Question 2: Selection Bias is positive or negative in the case?
- This means two groups could be quite different in other dimensions: other things are not equal.
- Observed association is *neither necessary nor sufficient for causality*.

# Observed Association: College vs Non-College Wage Differentials:

- Missing data problem also arises when we estimate ATE

i	$Y_{i1}$	$Y_{0i}$	$Y_i$	$D_i$	$Y_{i1} - Y_{0i}$
Tom	3	?	3	1	?
Jerry	2	?	2	1	?
Scarlett	?	1	1	0	?
Nicole	?	1	1	0	?
$E[Y_{1i} D_i = 1]$	$\frac{3+2}{2} = 2.5$				
$E[Y_{0i} D_i = 0]$	$\frac{1+1}{2} = 1$				
$E[Y_{1i} D_i = 1] - E[Y_{0i} D_i = 0]$	1.5				

- The Observed Association of attending college on average wage

$$\alpha_{corr} = 2.5 - 1 = 1.5$$

# Observed Association and Selection Bias

- But we are interested in causal effect, here is ATT

$$\alpha_{ATT} = E[\delta_i | D_i = 1] = E[Y_{1i} - Y_{0i} | D_i = 1] = 1$$

- So the selection bias

$$E[Y_{0i} | D_i = 1] - E[Y_{0i} | D_i = 0] = 0.5$$

- The **Selection Bias** is positive: *Those who attend college could be more intelligent so they can earn more even if they did not attend college.*

# Causal Effect and Identification Strategy

- Many Many Other examples
  - the effect of job training program on worker' s earnings
  - the effect of class size on students performance
  - ....
- Identification strategy tells us what we can learn about a causal effect from the available data.
- The main goal of identification strategy is **to eliminate the selection bias**.
- Identification depends on assumptions, not on estimation strategies.
- **“What's your identification strategy?”** = what are the assumptions that allow you to claim you' ve estimated a causal effect?

# Experimental Design as an Benchmark

# Randomized Controlled Trial

- A randomized controlled trial (RCT) is a form of investigation in which units of observation (e.g. individuals, households, schools, states) are randomly assigned to treatment and control groups.
- RCT has two features that can help us hold “other things equal” and then eliminates selection bias
  - Random assign treatment:
    - Randomly assign treatment (such as a coin flip) ensures that every observation has the same probability of being assigned to the treatment group.
    - Therefore, the probability of receiving treatment is unrelated to any other confounding factors.
  - Sufficient large sample
    - Large sample size can ensure that the group differences in individual characteristics wash out

# How to Solve the Selection Problem

- Random assignment of treatment  $D_i$  can eliminate selection bias. It means that the treated group is a random sample from the population.
- Being a random sample, we know that those included in the sample are **the same, on average**, as those not included in the sample on any measure.
- Mathematically, it makes  $D_i$  **independent** of potential outcomes, thus

$$D_i \perp (Y_{0i}, Y_{1i})$$

- **Independence:** Two variables are said to be independent if knowing the outcome of one provides no useful information about the outcome of the other.
  - Knowing outcome of  $D_i(0, 1)$  does not help us understand what potential outcomes of  $(Y_{0i}, Y_{1i})$  will be

# Random Assignment Solves the Selection Problem

- So we have

$$E[Y_{0i}|D_i = 1] = E[Y_{0i}|D_i = 0]$$

- Thus the **Selection Bias** equals to **ZERO**.
- Then **ATT** equals **Observed Association** because the

$$\begin{aligned} E[Y_{1i}|D_i = 1] - E[Y_{0i}|D_i = 0] &= E[Y_{1i}|D_i = 1] - E[Y_{0i}|D_i = 1] \\ &= E[Y_{1i} - Y_{0i}|D_i = 1] \end{aligned}$$

- No matter what assumptions we make about the distribution of  $Y$ , we can always estimate it with the difference in means.



# Our Benchmark: Randomized Experimental Methods

- Think of causal effects in terms of comparing counterfactuals or potential outcomes. However, we can never observe both counterfactuals —fundamental problem of causal inference.
- To construct the counterfactuals, we could use two broad categories of empirical strategies.
  - **Random Controlled Trials/Experiments:**
    - it can eliminates selection bias which is the most important bias arises in empirical research. If we could observe the counterfactual directly, then there is no evaluation problem, just **simply difference**.

# Our Benchmark: Randomized Experimental Methods

- We can generate the data of our interest by controlling experiments just as physical scientists or biologists do. But too obviously, we face more difficult and controversy situation than those in any other sciences.
- The various approaches using naturally-occurring data provide alternative methods of constructing the proper counterfactual
  - **Econometrics**
    - Congratuation! We are working and studying in a more tough and intractable area than others including most science knowledge.
- We should take the randomized experimental methods as our benchmark when we do empirical research whatever the methods we apply.

# Program Evaluation Econometrics

# Randomized Controlled Trials(RCT)

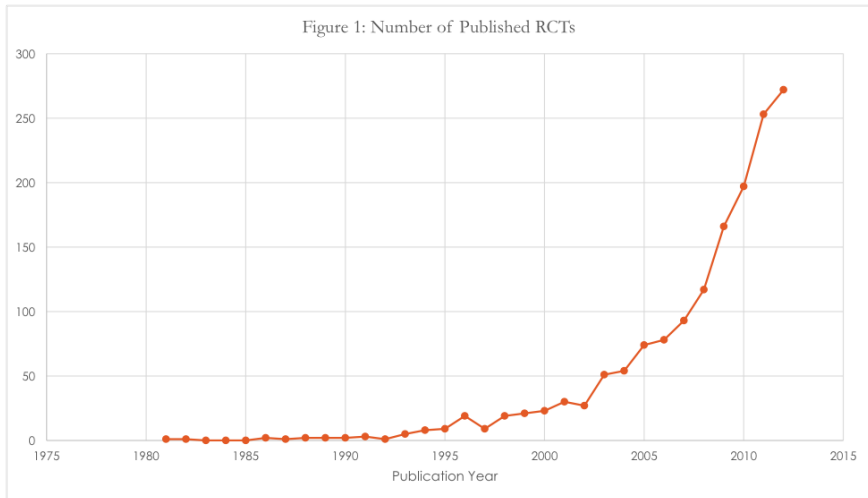
- First recorded RCT was done in 1747 by James Lind, who was a Scottish physician in the Royal Navy.
- Scurvy is a terrible disease caused by Vitamin C deficiency. Serious issue during long sea voyages.
- Lind took 12 sailors with scurvy and split them into six groups of two.
- Groups were assigned:
  - (1) 1 qt cider(苹果酒) (2) 25 drops of vitriol(硫酸) (3) 6 spoonfuls of vinegar, (4) 1/2 pt of sea water, (5) garlic, mustard(芥末) and barley water (大麦汤), (6) 2 oranges and 1 lemon
- Only Group 6 (citrus fruit) showed substantial improvement.

# Types of RCT

- Lab Experiments
  - eg: computer game for gamble in Lab
- Field Experiments
  - eg: the role of women in household' s decision or fake resumes in job application
- Quasi-Experiment or Natural Experiments: some unexpected institutional change or natural shock
  - eg: Germany reunion, Great Famine in China and U.S Bombing in Vietnam.

# Experiments and Publications

Figure:



# RCT are far from perfect!

- High Costs, Long Duration
- Potential Ethical Problems: “Parachutes reduce the risk of injury after gravitational challenge, but their effectiveness has not been proved with randomized controlled trials.”
  - Milgram Experiment
  - Stanford Prison Experiment
  - Monkey Experiment
- Limited generalizability
- RCTs allow us to gain knowledge about causal effects **without knowing the mechanism.**

# Potential Problems in Practice

- Small sample: Student Effect
- Hawthorne effect: The subjects are in an experiment can change their behavior.
- Attrition (样本流失): It refers to subjects dropping out of the study after being randomly assigned to the treatment or control group.
- Failure to randomize or failure to follow treatment protocol: People don't always do what they are told.
  - Wearing glasses program in Western Rural China.



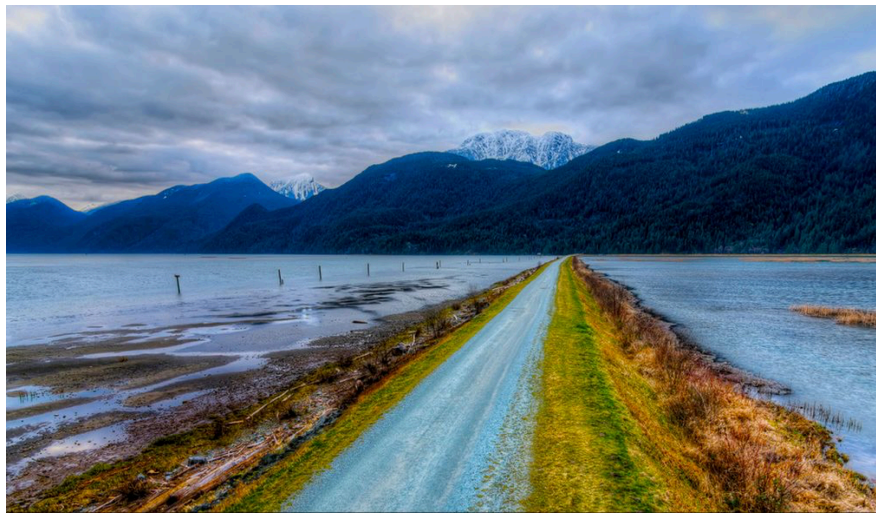
# Program Evaluation Econometrics(项目评估计量经济学)

- Question: How to do empirical research scientifically when we can not do experiments? It means that we always have selection bias in our data, or in term of “endogeneity” .
- Answer: Build a reasonable counterfactual world by naturally occurring data to find a proper control group is the core of econometrical methods.
- Here you **Furious Seven Weapons** in Applied Econometrics(七种盖世武器)
  - ① Random Controlled Trials(RCT)
  - ② OLS(回归)
  - ③ Matching (匹配)
  - ④ Decomposition (分解)
  - ⑤ Instrumental Variable (工具变量)
  - ⑥ Regression Discontinuity (断点回归)
  - ⑦ Differences in Differences (双差分) and Synthetic Control (合成控制法)

# Program Evaluation Econometrics(项目评估计量经济学)

- These Furious Seven are the most basic and popular methods in applied econometrics and so powerful that
  - even if you just master one, you may finish your empirical paper and get a good score.
  - if you master several ones, you could have opportunity to publish your paper.
  - If you master all of them, you might to teach applied econometrics class just as what I am doing now.
- We will introduce the essentials of these methods in the class as many as possible. Let's start our journey together.

# An Amazing Journey



# Let's Start Our Journey