

# 1. Potential Outcomes and Randomization

PhD Applied Methods

Duncan Webb  
NovaSBE

Spring 2026

# Intros

- **Me:** Duncan Webb
  - Development economist (field experiments in India, Madagascar, Colombia)
  - Email: dmbwebb@gmail.com – Office: B115B
  - Office hours: Wednesdays 1–2pm
- **You:** Name, research interests, what you hope to get from this course, what you did before starting the PhD

## Course Structure

- **Goal:** Deep understanding of modern causal inference methods
- **Format:** One 3-hour class per week
- **Topics:**
  - ① Potential Outcomes and Randomization (today)
  - ② Instrumental Variables
  - ③ Difference-in-Differences
  - ④ Regression Discontinuity Design
  - ⑤ Empirical Tools
- **Approach:** Theory + applications + implementation

## Assessment

- **Problem sets (20%):** 4 problem sets throughout the course
  - Applied/empirical questions (Stata, R, or Python)
  - Mathematical/theoretical questions
  - Individual work
  - Due Tuesdays at 12pm, starting next week
  - Submit to `dmbwebb@gmail.com`: **.tex**, **.pdf**, and **code files**
  - Include your name in the output
- **Final exam (80%):** 1.5-hour closed-book exam
  - Theoretical derivations and proofs
  - Applied reasoning and interpretation

## Assessment

- **AI policy:** You may use AI assistants, but 80% of your grade is a closed-book exam – AI only helps if you actually learn
- **Course materials:**  
<https://github.com/dmbwebb/NovaSBE-PhD-Econometrics-Students>
- **Questions?**

# Causality and understanding the world

- “We do not have knowledge of a thing until we have grasped its why, that is to say, its cause.” ~ Aristotle
- Not all research estimates a causal relationship, but the implication or takeaway of a paper is **almost always** a causal one
- Particularly important for **policy evaluation**:
  - What is the effect of microfinance on consumption?
  - What is the effect of reducing class size on education outcomes?
  - What is the effect of reducing the price of a good on its consumption?

## Causality and correlation

- For a very long time, economists made causal claims based on **correlations** and very shaky assumptions
- Until the **credibility revolution** (Angrist and Pischke, 2010) which formalized the conditions under which we could claim causality
  - And the use of **randomization** (or quasi-random events) to make those claims
- **Goal for this class:** Deep understanding of the tools we can use to make causal claims

## Rubin's potential outcomes framework

- The **potential outcomes framework** gives us a precise framework for thinking about when we can correctly claim to estimate the causal effect of some treatment
- The goal is to estimate the **causal effect** of some treatment, e.g.,
  - “Small class at school”
  - “Job training program”
  - ...
- Note that you can also have multiple treatments (e.g., small class, medium class, large class) and continuous treatments (University fees), but we’ll get to that later



# Counterfactuals

- **Outcome measure** is the outcome we care about
- Let  $Y_i$  be the observed outcome for individual  $i$
- For example:
  - $i$ 's wages in adulthood when examining impact of a job training program
  - $i$ 's test scores at school (for class size)
  - How much  $i$  discriminates against a minority (for prejudice-reduction intervention)



# Counterfactuals

- If  $i$  is not treated ( $D_i = 0$ ) then we **only** observe  $Y_i(0)$  and  $Y_i(1)$  is an unobserved counterfactual
- If  $i$  is treated ( $D_i = 1$ ) then we **only** observe  $Y_i(1)$  and  $Y_i(0)$  is an unobserved counterfactual

## Counterfactual quiz

Let's say our "treatment" ( $D_i$ ) is a **job training program**:

- What is the observed counterfactual for someone in the job training?
- What is  $Y_i(0)$  for someone in the control group?
- What is  $Y_i(0)$  for someone in the job training?
- Can we observe  $Y_i(1)$  for someone who doesn't get the training?
- What are  $Y_i(0)$  and  $Y_i(1)$  for someone who isn't even in the data?

## Causal effects

Using this framework, how would we write the **causal effect of the treatment on individual  $i$** ?

$$\Delta_j := Y_j(1) - Y_j(0) \quad (1)$$

This is the main thing we are trying to estimate!

In general,  $Y_i(1)$  and  $Y_i(0)$  can be different across people, and so  $\Delta_i$  may be different for each person too (“heterogeneous treatment effects”)





## Heterogeneous treatment effects

Call  $D_i = 0$  if untreated and  $D_i = 1$  if treated

Because the effect can be heterogenous, many evaluation parameters. In particular:

$$ATE := \mathbb{E}[Y_i(1) - Y_i(0)] \quad (2)$$

$$ATT := \mathbb{E}[Y_i(1) - Y_i(0)|D_i = 1] \quad (3)$$

- **ATE:** Average treatment effect - all population
- **ATT:** Average treatment on the treated - treated only (for instance, weak students are treated first)



## Observed outcomes

How do we compactly write the **actually observed outcome**?

$$Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0) \quad (4)$$

If  $D_i = 1$  then we observe  $Y_i(1)$

If  $D_i = 0$  then we observe  $Y_i(0)$

Think of it as a “binary switch”

NB: We can equivalently write this as  $Y_i = Y_i(0) + D_i(Y_i(1) - Y_i(0)) = Y_i(0) + D_i\Delta_i$  (i.e., in terms of the effect)

## Second quiz

- What is the observed outcome when  $D_i = 1$ ?
- What is the unobserved counterfactual when  $D_i = 0$ ?

## Regression with constant treatment effects

Consider a simple regression model:

$$Y_i = \alpha + D_i\beta + u_i \quad (5)$$

**Question:** What counterfactuals  $Y_i(0)$  and  $Y_i(1)$  generate this model?

## What does OLS estimate?

What does the OLS estimator of  $\beta$  measure in this regression?

$$Y_i = \alpha + D_i\beta + u_i \quad (6)$$

## OLS estimates the difference in means

Under  $\mathbb{E}[u_i|D_i] = 0$  (no selection bias):

$$\mathbb{E}[Y_i|D_i = 1] - \mathbb{E}[Y_i|D_i = 0] = \alpha + \beta - \alpha + \mathbb{E}[u_i|D_i = 1] - \mathbb{E}[u_i|D_i = 0] = \beta \quad (7)$$

So OLS gives us  $\beta$ . But what exactly is  $\beta$  in terms of treatment effects?

(Note: This model assumes constant treatment effects, i.e.,  $\Delta_i = \beta$  for all  $i$ . We'll revisit what happens when effects are heterogeneous later.)

## Outline

## What do we observe?

Let's think again about what can we actually observe in the data?

$$\mathbb{E}[Y_i(0)|D_i = 0]? \quad (8)$$

$$\mathbb{E}[Y_i(1)|D_i = 1]? \quad (9)$$

$$\mathbb{E}[Y_i(1)|D_i = 0]? \quad (10)$$

$$\mathbb{E}[Y_i(0)|D_i = 1]? \quad (11)$$

## What do we observe?

Identification problem: we “**observe**”

$$\mathbb{E}[Y_i(0)|D_i = 0] \quad (12)$$

$$\mathbb{E}[Y_i(1)|D_i = 1] \quad (13)$$

but not the counterfactuals

$$\mathbb{E}[Y_i(1)|D_i = 0] \quad (14)$$

$$\mathbb{E}[Y_i(0)|D_i = 1] \quad (15)$$

But what would we need to estimate the  $ATT$ , for instance?

$$ATT = \overbrace{\mathbb{E}[Y_i(1)|D_i = 1]}^{\text{"observed"}} - \overbrace{\mathbb{E}[Y_i(0)|D_i = 1]}^{\text{"unobserved"}} \quad (16)$$



**Hypothesis to identify *ATT*:**

$$\mathbb{E}[Y_i(0)|D_i = 1] = \mathbb{E}[Y_i(0)|D_i = 0] = \mathbb{E}[Y_i(0)]$$

i.e. no selectivity: treated “are like” untreated

Then

$$ATT = \overbrace{\mathbb{E}[Y_i(1)|D_i = 1]}^{\text{"observed"}} - \overbrace{\mathbb{E}[Y_i(0)|D_i = 1]}^{\text{"unobserved"}} \quad (17)$$

$$= \overbrace{\mathbb{E}[Y_i(1)|D_i = 1]}^{\text{“observed”}} - \overbrace{\mathbb{E}[Y_i(0)|D_i = 0]}^{\text{“observed”}} \quad (18)$$

In this very simple case, compare empirical means in each group

The counterfactual for a group is simply the observed outcome of the other group

NB: to identify  $ATE$  we need more:

$$\mathbb{E}[Y_i(0)|D_i] = \mathbb{E}[Y_i(0)] \quad \text{and} \quad \mathbb{E}[Y_i(1)|D_i] = \mathbb{E}[Y_i(1)]$$

$$ATE = \overbrace{\mathbb{E}[Y_i(1)]}^{\text{"unobserved"}} - \overbrace{\mathbb{E}[Y_i(0)]}^{\text{"unobserved"}} \quad (19)$$

$$= \overbrace{\mathbb{E}[Y_i(1)|D_i = 1]}^{\text{"observed"}} - \overbrace{\mathbb{E}[Y_i(0)|D_i = 0]}^{\text{"observed"}} \quad (20)$$

Under these assumptions, we also identify the *ATT* - **why?**

⇒ Because we assumed that the counterfactuals are similarly distributed (on average) in both populations



## Why is selection bias quite likely?

Simple **Roy model**: “I am in if this is worth it”

$$D_i = 1 \text{ if } Y_i(1) - Y_i(0) > c$$

Then, in general

$$\mathbb{E}[Y_i(0)|D_i = 1] = \mathbb{E}[Y_i(0)|Y_i(0) < Y_i(1) - c] \quad (24)$$

$$\neq \mathbb{E}[Y_i(0)|Y_i(0) \geq Y_i(1) - c] = \mathbb{E}[Y_i(0)|D_i = 0] \quad (25)$$

In this case, selectivity stems from

- **Comparative advantages** ( $Y_i(1) - Y_i(0)$  large for some, small for others). Most simple instance: participants have smaller  $Y_i(0)$ , thus larger potential gain
- **Heterogeneity in cost  $c$**  (if it is correlated with  $Y_i(0)$ )

## Why is selection bias quite likely?

### Another reason for selection bias: **administrative rules**

For instance:

- **“Cream-skimming”**: they choose “the best”, and  $\mathbb{E}[Y_i(0)|D_i = 1] > \mathbb{E}[Y_i(0)|D_i = 0]$
- **Remedial targeting**: e.g. focus on intervening with weak kids in the class, so  $\mathbb{E}[Y_i(0)|D_i = 1] < \mathbb{E}[Y_i(0)|D_i = 0]$

## Link with endogeneity

**Selectivity** will lead to **endogeneity** of  $D_i$  in a regression

For simplicity, focus on a simple model with homogenous effects, i.e.  $\Delta_i = \beta$  for everyone, and  $u_1 = u_0 = u_i$ :

$$Y_i = \alpha + D_i\beta + u$$

Selectivity is then:

$$\mathbb{E}[Y_i(0)|D_i = 1] \neq \mathbb{E}[Y_i(0)|D_i = 0] \quad (26)$$

$$\Rightarrow \mathbb{E}[u_j | D_j = 1] \neq \mathbb{E}[u_j | D_j = 0] \quad (27)$$



## A second problem: Heterogeneous treatment effects

We've seen that **selection bias** is a major obstacle to causal inference.

But there's a **second problem** we need to consider:

## What if treatment effects **vary across individuals**?

That is, what if  $\Delta_i = Y_i(1) - Y_i(0)$  differs from person to person?

This matters because even **without selection bias**, OLS may not estimate ATE or ATT if effects are heterogeneous.



## Are treatment effects constant?

**Question:** Is it realistic that  $\Delta_i = Y_i(1) - Y_i(0)$  is the same for everyone?

### Examples where effects vary:

- Job training: More effective for workers with less experience
- Class size reduction: May help struggling students more
- Medicine: Effects vary with age, weight, genetics
- Education subsidies: Returns higher for high-ability students

⇒ In most applications, treatment effects are **heterogeneous** across individuals

## A more general model

Allow potential outcomes to differ flexibly across individuals:

$$Y_i(0) = g_0(X_i) + u_{0i} \quad (32)$$

$$Y_i(1) = g_1(X_i) + u_{1i} \quad (33)$$

where:

- $g_0, g_1$  = functions of observable characteristics  $X_i$
- $u_{0i}, u_{1i}$  = unobservable components (can differ by treatment status)

**Key difference from simple model:**

- The effect can vary with  $X_i$  (observable heterogeneity)
- The effect can vary with  $u_{1i} - u_{0i}$  (unobservable heterogeneity)

## Heterogeneous treatment effects

Recall:  $Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0)$

Substituting our general model:

$$Y_i = g_0(X_i) + D_i \underbrace{[g_1(X_i) - g_0(X_i) + u_{1i} - u_{0i}]}_{\Delta_i} + u_{0i} \quad (34)$$

The individual treatment effect is now:

$$\Delta_i = g_1(X_i) - g_0(X_i) + u_{1i} - u_{0i} \quad (35)$$

This varies across individuals—it's a **random coefficient!**

## Does OLS estimate ATE or ATT?

**Question:** With heterogeneous effects, does OLS estimate ATE or ATT?

OLS estimates the difference in means:

$$\mathbb{E}[Y_i|D_i = 1] - \mathbb{E}[Y_i|D_i = 0] = \mathbb{E}[g_1(X_i) + u_{1i}|D_i = 1] - \mathbb{E}[g_0(X_i) + u_{0i}|D_i = 0] \quad (36)$$

**Problem:** OLS compares  $g_1$  for *treated* against  $g_0$  for *untreated*

But treated and untreated may have different  $X_i$  distributions.

**Key result:** Even if  $\mathbb{E}[u_{ki}|D_i] = 0$  (no selection bias), OLS estimates **neither** ATE nor ATT in general.

## Summing up: Two problems for causal inference

### Problem 1: Selection bias

- Treated and untreated differ in ways that affect outcomes
- $\mathbb{E}[Y_i(0)|D_i = 1] \neq \mathbb{E}[Y_i(0)|D_i = 0]$
- OLS confounds treatment effect with pre-existing differences

## Problem 2: Heterogeneous effects

- Treatment effects vary across individuals ( $\Delta_i$  not constant)
- Even without selection bias, OLS may not estimate ATE or ATT
- Depends on how  $X_i$  is distributed across treatment groups

**Question:** How can we solve **both** problems?

## Outline

1. Potential outcomes framework
2. Selection bias
3. Controlled experiments
4. Conclusion

## Idea behind experiments

Simplest way to identify treatment causal effect: make likely the hypotheses

$$\mathbb{E}[Y_i(\mathbf{0})|D_i = 1] = \mathbb{E}[Y_i(\mathbf{0})|D_i = 0] = \mathbb{E}[Y_i(\mathbf{0})] \quad (37)$$

$$\mathbb{E}[Y_i(\mathbf{1})|D_i = 1] = \mathbb{E}[Y_i(\mathbf{1})|D_i = 0] = \mathbb{E}[Y_i(\mathbf{1})] \quad (38)$$

If we draw treated and untreated **randomly** from a population then:

$$ATT = \overbrace{\mathbb{E}[Y_i(1)|D_i = 1]}^{\text{"observed"}} - \overbrace{\mathbb{E}[Y_i(0)|D_i = 1]}^{\text{"unobserved"}} \quad (39)$$

$$= \overbrace{\mathbb{E}[Y_i(1)|D_i = 1]}^{\text{"observed"}} - \overbrace{\mathbb{E}[Y_i(0)|D_i = 0]}^{\text{"observed"}} \quad (40)$$

Can be estimated with empirical means (or via regressions)

## Idea behind experiments

- **Intuition:** if we randomly select who receives the treatment and who doesn't, then **on average** it will be similar types of people in each group, and so the average counterfactuals will be the same
- Therefore, any difference we **do** observe after the treatment must be **caused by the treatment**
- This is why randomized controlled trials are called the **gold standard** of evidence (somewhat controversially)
- **Other methods for causal inference** are built on this paradigm – other identification methods “mimic” random assignment into treatment



## Randomization solves both problems

### Problem 1: Selection bias

- Randomization ensures  $\mathbb{E}[Y_i(0)|D_i = 1] = \mathbb{E}[Y_i(0)|D_i = 0]$
- No systematic differences between treated and untreated groups

## Problem 2: Heterogeneous effects

- Randomization ensures  $X_i$  is distributed identically across groups
- So  $\mathbb{E}[g_k(X_i)|D_i = 1] = \mathbb{E}[g_k(X_i)|D_i = 0]$  for  $k = 0, 1$

**Key insight:** Randomization ensures  $OLS = ATE$ , even with heterogeneous treatment effects!

## Use of randomized controlled trials

- **Popularity:** AER+JPE+QJE, 0.8% of published articles in 1983 → 8.2% in 2011 (while theory: 58% → 19%)
- **Nobel Prize** in economics to Esther Duflo, Abhijit Banerjee, Michael Kremer for pioneering this methodology in development economics.
- **Infrastructure** - organisations like J-PAL and IPA provide infrastructure for this kind of research

## Critiques of randomized controlled trials

- **Equilibrium effects** - standard methodology ignores spillover effects or equilibrium effects, although frontier methods and large trials can understand these (see e.g., Egger et al, *Econometrica* 2022)
- **External validity** - how to generalize from the context of the RCT to other contexts
- **Ethics** - is it ethical to deny treatment to the control group? This depends on the context, and what the alternative is, e.g., it's no longer ethical to deny proven medical treatment
- **Mechanisms** - early RCTs tried to measure treatment effects, but high-quality studies now focus a lot on understanding mechanisms using additional treatments or by explicitly testing theory-driven models

## Simplest design

Randomize individuals and compare treated and untreated:

$$\bar{Y}_1 - \bar{Y}_0 = \frac{1}{N_1} \sum_{i \in D_1=1} Y_i - \frac{1}{N_0} \sum_{i \in D_1=0} Y_i$$

Similar to OLS on

$$y = \alpha + \beta D_i + u$$

(Leave as an exercise to prove it algebraically using OLS matrix formula)

## Real life design issues in RCTs

- 1 **Balance** between treatment groups in a finite sample
- 2 **Adding controls**
- 3 **Imperfect compliance**

## Real example: reducing class size in “STAR” program

## EXPERIMENTAL ESTIMATES OF EDUCATION PRODUCTION FUNCTIONS\*

ALAN B. KRUEGER

This paper analyzes data on 11,600 students and their teachers who were randomly assigned to different size classes from kindergarten through third grade. Statistical methods are used to adjust for nonrandom attrition and transitions between classes. The main conclusions are (1) on average, performance on standardized tests increases by four percentile points the first year students attend small classes; (2) the test score advantage of students in small classes expands by about one percentile point per year in subsequent years; (3) teacher aides and measured teacher characteristics have little effect; (4) class size has a larger effect for minority students and those on free lunch; (5) *Hawthorne* effects were unlikely.

## 1. Balance between treatment groups

TABLE I  
COMPARISON OF MEAN CHARACTERISTICS OF TREATMENTS AND CONTROLS:  
UNADJUSTED DATA

A. Students who entered STAR in kindergarten <sup>b</sup>				
Variable	Small	Regular	Regular/Aide	Joint <i>P</i> -Value <sup>a</sup>
1. Free lunch <sup>c</sup>	.47	.48	.50	.09
2. White/Asian	.68	.67	.66	.26
3. Age in 1985	5.44	5.43	5.42	.32
4. Attrition rate <sup>d</sup>	.49	.52	.53	.02
5. Class size in kindergarten	15.1	22.4	22.8	.00
6. Percentile score in kindergarten	54.7	49.9	50.0	.00

## 2. Adding controls with an OLS regression of an RCT

This shows the effect of  $D_i$  ("small class") on  $Y_i$  (percentile on standardized test score).

TABLE V  
OLS AND REDUCED-FORM ESTIMATES OF EFFECT OF CLASS-SIZE ASSIGNMENT ON  
AVERAGE PERCENTILE OF STANFORD ACHIEVEMENT TEST

Explanatory variable	Reduced form: initial class size			
	(5)	(6)	(7)	(8)
Small class	4.82 (2.19)	5.37 (1.25)	5.36 (1.21)	5.37 (1.19)
Regular/aide class	.12 (2.23)	.29 (1.13)	.53 (1.09)	.31 (1.07)
White/Asian (1 = yes)	—	—	8.35 (1.35)	8.44 (1.36)
Girl (1 = yes)	—	—	4.48 (.63)	4.39 (.63)
Free lunch (1 = yes)	—	—	-13.15 (.77)	-13.07 (.77)
White teacher	—	—	—	-.57 (2.10)
Teacher experience	—	—	—	.26 (.10)
Master's degree	—	—	—	-.51 (1.06)
School fixed effects	No	Yes	Yes	Yes
$R^2$	.01	.25	.31	.31



## Adding controls

## OLS often used because it allows you to add controls: **Why?**

If assignment is truly random, conditioning on  $X_i$  should not affect point estimates

We have  $\mathbb{E}[Y_i(0)|D_i] = \mathbb{E}[Y_i(0)]$  and  $\mathbb{E}[X_i|D_i] = \mathbb{E}[X_i]$

Therefore OLS on  $Y_i = X_i\gamma + \beta D_i + u$  gives (asymptotically) the same  $\beta$  as  $Y_i = \beta D_i + u'_i$  (Frisch-Waugh theorem)

## Adding controls

**But** it's still useful because it increases precision: **Why?**

$$V(\hat{\beta}) = \sigma^2(X'X)^{-1}$$

where  $\sigma^2$  is residual variance

$$y = \beta D_j + u' \quad (41)$$

$$y = \beta D_i + x\gamma + u \quad (42)$$

$$V(u') = V(x\gamma) + V(u) > V(u)$$

Thus, the second equation estimates the same  $\beta$  but with more precision

Depends on how much  $X_i$  explain  $Y_i$  (and may not hold in finite samples)

## Controls and precision

Effect of  $D_i$  ("small class") on  $Y_i$  (percentile on standardized test score):

TABLE V  
OLS AND REDUCED-FORM ESTIMATES OF EFFECT OF CLASS-SIZE ASSIGNMENT ON  
AVERAGE PERCENTILE OF STANFORD ACHIEVEMENT TEST

Explanatory variable	Reduced form: initial class size			
	(5)	(6)	(7)	(8)
Small class	4.82 (2.19)	5.37 (1.25)	5.36 (1.21)	5.37 (1.19)
Regular/aide class	.12 (2.23)	.29 (1.13)	.53 (1.09)	.31 (1.07)
White/Asian (1 = yes)	—	—	8.35 (1.35)	8.44 (1.36)
Girl (1 = yes)	—	—	4.48 (.63)	4.39 (.63)
Free lunch (1 = yes)	—	—	-13.15 (.77)	-13.07 (.77)
White teacher	—	—	—	-.57 (2.10)
Teacher experience	—	—	—	.26 (.10)
Master's degree	—	—	—	-.51 (1.06)
School fixed effects	No	Yes	Yes	Yes
$R^2$	.01	.25	.31	.31

## Stratified randomization (block randomization)

**Motivation:** With simple randomization, we might get imbalances on important characteristics (especially in small samples)

**Simple randomization:** Randomly assign all units to treatment or control

**Stratified randomization:** Divide sample into strata based on pre-treatment characteristics, then randomize **within each stratum**

**Example:** Stratify by education (e.g., high school vs college) and gender  $\Rightarrow$  4 strata. Within each, randomly assign 50% to treatment.

**Key difference:** Stratified randomization ensures representation from each subgroup

# Why use stratified randomization?

Two main benefits:

## 1. Ensures balance on stratification variables

- With simple randomization, treatment and control groups may differ on key characteristics (especially with small samples)
- Stratification **guarantees** balance on the stratification variables
- Example: Exactly 50% of treated are female if you stratify by gender

## 2. Increases precision (lowers standard errors)

- If stratification variables predict outcomes, controlling for strata reduces residual variance
- Recall:  $V(\hat{\beta}) = \sigma^2(X'X)^{-1}$  where  $\sigma^2$  is residual variance
- Lower residual variance  $\Rightarrow$  smaller standard errors  $\Rightarrow$  more statistical power

# Implementing stratified randomization

## How to choose stratification variables?

- Pick variables that are:
  - Strong predictors of the outcome (increases precision)
  - Measured before randomization (ensures exogeneity)
  - Create a manageable number of strata (rule of thumb: at least 4-6 observations per stratum-treatment combination)
- Common choices: baseline outcome, gender, age groups, geographic location

# Implementing stratified randomization

## Specification with stratified randomization:

Include **stratum fixed effects** in your regression:

$$Y_i = \alpha + \beta D_i + \sum_{s=1}^S \gamma_s \mathbb{1}[\text{Stratum}_i = s] + u_i \quad (43)$$

- This accounts for how randomization was done
- Improves precision (even though  $\beta$  estimate is similar without FE)
- Standard practice: always control for strata used in randomization

### 3. Imperfect Compliance

So far we've assumed people **comply** with their assignment.

But what if they don't?

- What happens when some people assigned to treatment don't take it?
- What happens when some people assigned to control get treated anyway?
- Can we still estimate causal effects? If so, *what* causal effects?



## Imperfect Compliance: The Problem

In practice, we often **cannot force people to comply** with their assignment

Two types of non-compliance:

- **Non-take-up:** Assigned to treatment but don't take it
- **Crossover:** Assigned to control but take treatment anyway

**Key question:** How does this affect our ability to estimate causal effects?

## Example: Krueger Class Size Experiment

### What went wrong?

- ① Approx. 10% changed class type during the experiment
  - Teacher requests (behavioral problems)
  - Parent pressure
- ② Some children changed school or moved (“attrition”)

This is called an **encouragement design**:

- We *encourage* but don’t *force* treatment
- Simpler to implement, more acceptable, often no choice
- Comes at a cost to precision

**Question:** What can we still learn?

# The Fundamental Insight

**Problem:** We can no longer directly compare treated vs untreated

- Actual treatment receipt ( $T_i$ ) is now *endogenous*—a choice!
- Selection bias is back

**But:** We can still compare *assigned* vs *not assigned*

- Assignment ( $D_i$ ) is still *random*
- This comparison is “clean”

This gives us **two distinct questions**:

- ① What is the effect of being *assigned* to treatment? (ITT)
- ② What is the effect of actually *receiving* treatment? (Wald/IV)

# Intention-to-Treat (ITT)

## Intention-to-Treat (ITT):

$$ITT = \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0]$$

The ITT compares outcomes by **assignment**, ignoring actual treatment

**Is this causal?** Yes! Assignment is random

**Interpretation:** The effect of being assigned/offered/encouraged to treatment

## When is ITT Useful?

**Policy perspective:** If you're implementing a program, ITT tells you the effect of *your policy* (including non-compliance)

**Example:** Vaccine rollout

- ITT captures real-world effectiveness
- Accounts for people who won't show up, refuse, etc.
- This is what a policymaker cares about!

**Limitation:** Doesn't tell us the effect of the treatment *itself*

**Question:** What if we want to know the effect of actually receiving treatment?

## Notation: Assignment vs Treatment

We need to distinguish two things:

- $D_i \in \{0, 1\}$ : Random **assignment** (what we control)
- $T_i \in \{0, 1\}$ : Actual **treatment received** (what we observe)

Define compliance probabilities:

$$p_1 = P(T_i = 1 | D_i = 1) \quad - \text{compliance rate among assigned-to-treatment} \quad (44)$$

$$p_0 = P(T_i = 1 | D_i = 0) \quad - \text{crossover rate among assigned-to-control} \quad (45)$$

**Perfect compliance:**  $p_1 = 1$  and  $p_0 = 0$

**Imperfect compliance:**  $p_1 < 1$  or  $p_0 > 0$  (or both)

# Visualizing Compliance

$D = 0$ <i>Assigned to control</i>	$D = 1$ <i>Assigned to treatment</i>
<p><b>Stay in control:</b>  <math>(1 - p_0)</math> fraction  <math>T_i = 0</math>            -----</p> <p><b>Cross to treatment:</b>  <math>p_0</math> fraction  <b>TREATED (<math>T_i = 1</math>)</b></p>	<p><b>TREATED (<math>T_i = 1</math>)</b>  <math>p_1</math> fraction</p>

**Key insight:** Only  $D$  groups are comparable (random), not  $T$  groups!

The choice to cross over is *endogenous*

# From ITT to Treatment Effect: Intuition

**Core logic:**

$$\text{ITT} = (\text{effect of treatment}) \times (\text{change in treatment probability})$$

↓

$$\text{Treatment effect} = \frac{\text{ITT}}{\text{change in treatment probability}}$$

**Intuition:** If assignment shifts treatment probability by 50%, and outcomes improve by 3, then treatment must improve outcomes by  $3/0.5 = 6$



## Numerical Example Setup

**Setting:** 8 students, assigned to small ( $D = 1$ ) or large ( $D = 0$ ) class

$D = 0$ <i>Assigned to large class</i>	$D = 1$ <i>Assigned to small class</i>
1: untreated – score 5	5: treated – score 17
2: untreated – score 5	6: treated – score 5
3: treated – score 15	7: treated – score 15
4: treated – score 15	8: treated – score 15

**Compliance:**

- All  $D = 1$  students get treatment:  $p_1 = 1$
- 2 of 4  $D = 0$  students *also* get treatment:  $p_0 = 0.5$

# Calculating ITT and Treatment Effect

**Step 1:** Calculate mean outcomes by assignment

- $\mathbb{E}[Y|D = 1] = (17 + 5 + 15 + 15)/4 = 13$
- $\mathbb{E}[Y|D = 0] = (5 + 5 + 15 + 15)/4 = 10$

**Step 2:** ITT

$$ITT = 13 - 10 = 3$$

**Step 3:** Treatment effect

$$\tau = \frac{ITT}{p_1 - p_0} = \frac{3}{1 - 0.5} = \frac{3}{0.5} = 6$$

**Interpretation:** Moving a student from large to small class increases score by 6 points

# The Wald Estimator: Formal Derivation

Under constant treatment effects, we can derive:

$$\mathbb{E}[Y|D = 1] = \mathbb{E}[Y_i(0)] + \tau \cdot p_1 \quad (46)$$

$$\mathbb{E}[Y|D = 0] = \mathbb{E}[Y_i(0)] + \tau \cdot p_0 \quad (47)$$

Subtracting:

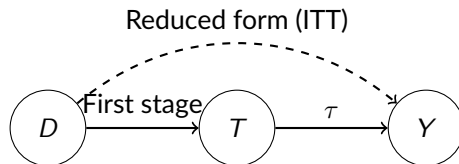
$$\mathbb{E}[Y|D = 1] - \mathbb{E}[Y|D = 0] = \tau(p_1 - p_0)$$

Solving for  $\tau$ :

$$\tau = \frac{\mathbb{E}[Y|D = 1] - \mathbb{E}[Y|D = 0]}{P(T = 1|D = 1) - P(T = 1|D = 0)} = \frac{\text{ITT}}{p_1 - p_0}$$

Both numerator and denominator are **observable**

## Reduced Form and First Stage



- **Reduced form** ( $D \rightarrow Y$ ): Effect of assignment on outcome = ITT
- **First stage** ( $D \rightarrow T$ ): Effect of assignment on treatment =  $p_1 - p_0$
- **Wald estimator** = Reduced form / First stage

## Two-Stage Least Squares (2SLS)

**Implementation** as a simultaneous equation model:

**First stage:** Predict treatment from assignment

$$\hat{T}_i = \hat{\pi}_0 + \hat{\pi}_1 D_i$$

**Second stage:** Regress outcome on *predicted* treatment

$$Y_i = \beta_0 + \tau \hat{T}_i + u_i$$

$D$  is an “**instrument**” for  $T$ :

- $D$  is exogenous (random assignment)
- $D$  affects  $Y$  only through  $T$  (exclusion restriction)
- $D$  predicts  $T$  (relevance:  $\pi_1 \neq 0$ )

This is **instrumental variables**—more detail next lecture!

## Important Caveat: Who Does This Apply To?

**Warning:** The Wald/IV estimate is a **Local Average Treatment Effect (LATE)**

It applies to “**compliers**”—those whose treatment status is *changed* by assignment

This is **NOT the same** as the ATE in general:

- People who always take treatment regardless of assignment? Not included
- People who never take treatment regardless of assignment? Not included
- Only those who *comply* with their assignment are captured

**Preview:** We’ll study this carefully in the IV lecture—understanding who the compliers are is crucial for interpreting IV estimates

## The Cost of Non-Compliance

## What is the cost of low compliance? Lower precision in our estimates

With **full compliance**, estimating  $Y_i = \alpha + \tau D_i + u_i$  gives:

$$V(\hat{\tau}) = \frac{1}{\bar{D}(1 - \bar{D})} \cdot \frac{V(u)}{N}$$

With **imperfect compliance**, the IV estimator has variance:

$$V(\hat{\tau}_{IV}) = \frac{1}{\bar{D}(1 - \bar{D})} \cdot \frac{V(u)}{N} \cdot \frac{1}{(p_1 - p_0)^2}$$

**Key insight:** Standard errors are inflated by factor  $\frac{1}{p_1 - p_0}$





## Intuition: Why Does Compliance Matter?

**Intuition:** With low compliance,  $D$  is a “blurry lever” for  $T$

- You **don't know** exactly who was induced to comply by the assignment
- Assignment is a weak predictor of actual treatment
- This makes it harder to detect the treatment effect signal

**With volunteers:** You know exactly who is complying

**Trade-off:** External validity (volunteers may differ from population) vs precision

## Summary: Two Approaches to Imperfect Compliance

	ITT	Wald/IV
<b>Estimand</b>	Effect of assignment	Effect of treatment
<b>Formula</b>	$\mathbb{E}[Y D = 1] - \mathbb{E}[Y D = 0]$	$\frac{\text{ITT}}{p_1 - p_0}$
<b>Identified?</b>	Always	Requires $p_1 \neq p_0$
<b>Applies to</b>	Everyone	Compliers only
<b>Policy use</b>	Program effectiveness	Treatment efficacy

**Key insight:** The choice between ITT and IV depends on your research question

- Policy evaluation? → ITT tells you what your program achieves
- Treatment efficacy? → IV tells you what the treatment does (for compliers)

## Outline

## Summing up

- 1 **Potential outcomes framework** - this gives us a way of conceptualizing counterfactuals and articulating clearly when we can and cannot make causal claims
- 2 **Randomized controlled trials** are a way to make causal claims with relatively weak assumptions on the data generating process
- 3 **Design issues** - we learnt about various design issues that come up in RCTs, e.g., dealing with imbalances, adding controls, and dealing with imperfect compliance