

# Longitudinal Data Analysis

---

White Rose Social Sciences Doctoral Training Partnership

 Thiago Oliveira

 Lecturer in Criminology, University of Manchester

 8 April 2025, University of Leeds

## Part II: causal inference with panel data

---

# Potential outcomes framework

---

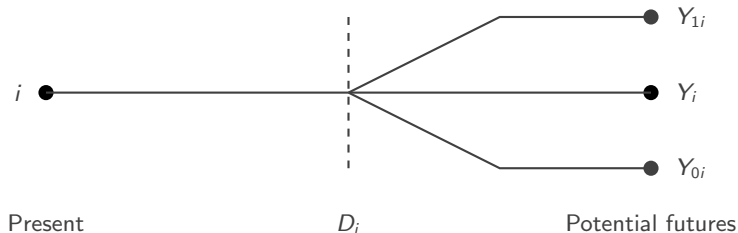
## Causal inference and the potential outcomes framework

---

Goal in causal inference is to assess the causal effect of a treatment/exposure on some outcome

- ↪ Does raising the minimum wage reduce employment?
- ↪ Does housing assistance reduce homelessness?
- ↪ Do body-worn cameras reduce police use of force?
- ↪ Does voting by mail increase voter turnout?
- ↪ Does exposure to misinformation reduce political trust??
- ↪ ...

## Causal inference and the potential outcomes framework



# Causal inference and the potential outcomes framework

---

$Y_i$ : Observed outcome variable of interest for unit  $i$

## Potential outcomes

$Y_{0i}$  and  $Y_{1i}$ : Potential outcomes for unit  $i$

$$Y_{\cdot i} = \begin{cases} Y_{1i} & \text{Potential outcome for unit } i \text{ with treatment} \\ Y_{0i} & \text{Potential outcome for unit } i \text{ without treatment} \end{cases}$$

$D_i$ : Indicator of treatment intake for *unit*  $i$

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

## Definition of causal effect

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

## Fundamental problem of causal inference

$\rightsquigarrow$  We cannot observe both potential outcomes for the same unit  $i$ !

# Causal inference and the potential outcomes framework

---

## Randomisation solves the problem!

Logic of randomised control trials

- ~> Randomly divide a sample in two groups
- ~> Because this was random, both groups are *on average* the same
- ~> Then apply the treatment/exposure to one group (the treatment group), but not the other (control group)
- ~> Because the exposure happened after the treatment assignment, the only difference between the two groups is the treatment/exposure
- ~> Therefore, any subsequently observed differences are attributable to the treatment/exposure
- ~> We randomisation, we can thus find the average treatment effect

## Causal inference and the potential outcomes framework

---

What if we cannot conduct an experiment?

↪ Randomised Experiments

↪ Observational Studies

- Selection on observables
  - Regression
  - Matching
  - Weighting
- Selection on unobservables
  - Difference-in-Differences and synthetic control
  - Instrumental Variables
  - Regression Discontinuity Designs



## Causal inference and the potential outcomes framework

---

- ↪ Causality is defined by potential outcomes, not by realised (observed) outcomes
  - ↪ Observed association is neither necessary nor sufficient for causality
  - ↪ Estimation of causal effects of a treatment (usually) starts with studying the assignment mechanism
  - ↪ The goal is to mimic the features of a randomised experiment even if we don't have one
  - ↪ When we don't have an RCT, our ability to make causal inferences often relies on making untestable assumptions about the assignment mechanism
- ⇒ Now let's see how we can leverage panel data to make causal inferences!

# Difference-in-differences

---

## Intuition of the difference-in-differences estimator

⇒ What if we use **time** in our favour?

- ↪ Collect data on  $Y$  at two points in time: **before** and **after** the treatment/exposure/policy intervention
- ↪ Analyse the extent to which  $Y$  **changes** in units that received the treatment
- ↪ Analyse the extent to which  $Y$  **changes** in units that did NOT receive the treatment
- ↪ Compare the two **changes**

## Intuition of the difference-in-differences estimator

Some conceptual clarification to make our lives easier

↪ Variation **between** units: difference

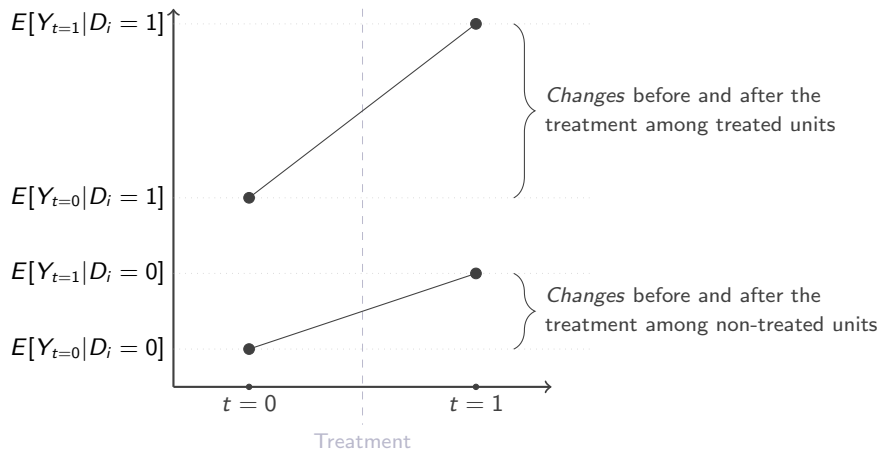
↪ Variation **within** units (over time): changes

⇒ We want to estimate the difference in changes  
or (*difference-in-differences*)

↪ The difference between (a) **changes in  $Y$  before and after the intervention among treated units** and (b) **changes in  $Y$  before and after the intervention among non-treated units** is the causal effect!

(under some assumptions regarding those changes... Let's dive into it)

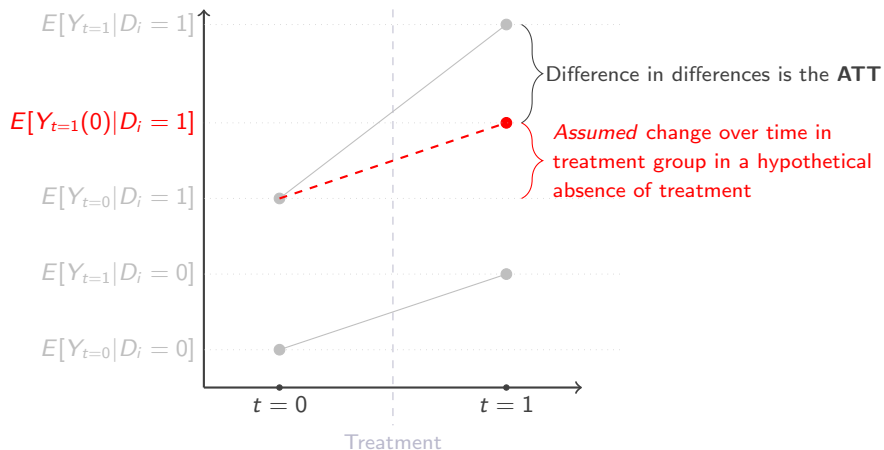
## Difference-in-differences setup



↪ **Problem:** Missing potential outcomes:  $E[Y_{i,t=1}(0) | D_i = 1]$  and  $E[Y_{i,t=1}(1) | D_i = 0]$

## Difference-in-differences setup

Strategy: Use the change in the control group to *assume*  $E[Y_{t=1}(0)|D_i = 1]$



Assumption: Trend over time would be the same for treatment and control

## Identification assumption

### Parallel trends

↪ Had the treated units not received the treatment, they would have followed the same trend as the control units

### Difference-in-differences estimator

Difference in changes:

$$\delta_{ATT} = \left\{ \begin{array}{l} \text{Changes in treatment group before and after treatment} \\ - \left\{ \text{Changes in control group before and after treatment} \right\} \end{array} \right\}$$

## Threats to validity

### Non-parallel trends

↪ **Very critical assumption:** treatment units have similar trends to control units in the absence of treatment

↪ **Fundamental problem of causal inference:** we cannot observe potential outcome under the control condition for treated units in the post-treatment period

⇒ **What can we do?**

*(more on that later...)*

- Careful assessment: is assuming parallel trends plausible?
- Estimate treatment effects at different time points (placebo tests)



# Using regression to estimate the difference-in-differences

---

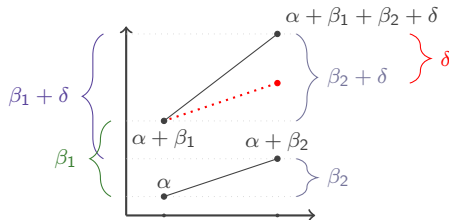
## Estimator (Regression 1)

We can obtain the difference in differences using regression techniques.

$$Y_i = \alpha + \beta_1 \cdot D_i + \beta_2 \cdot T_i + \delta \cdot (D_i \cdot T_i) + \varepsilon.$$

We can see that:

$E[Y_i   D_i, T_i]$	$T_i = 0$	$T_i = 1$	Changes after - before
$D_i = 0$	$\alpha$	$\alpha + \beta_2$	$\beta_2$
$D_i = 1$	$\alpha + \beta_1$	$\alpha + \beta_1 + \beta_2 + \delta$	$\beta_2 + \delta$
Treated - control	$\beta_1$	$\beta_1 + \delta$	$\delta$



# DiD: First differences estimator

## Estimator (Regression 2)

*With panel data we can use regression with first differences:*

$$\Delta Y_i = \alpha + \delta \cdot D_i + X' \beta + u,$$

*where  $\Delta Y_i = Y_i(1) - Y_i(0)$ .*

- *With two periods this gives the same result as other regressions*

# Advantages of the regression estimator

## 1. We can include covariates

- Controlling for some covariates may increase precision
- Time-varying covariates may strengthen the parallel assumptions
- (add covariates cautiously! e.g., beware of post-treatment bias)

## 2. Easy to calculate standard errors

- (though be careful about clustering)

## 3. Easy to extend to other types of treatment

- (not just binary)

## Some limitations

>> This setup only works for the simplest scenario with two time periods

↪ It doesn't make use more periods

- Useful to make careful assessments of time trends

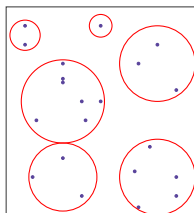
↪ Sometimes different units are treated at different time points

# Difference-in-differences with multiple periods

---

# Intuition of fixed-effect regression

>> Assume a pool of structured data



↪ Each **dot** represents a **unit  $i$**

↪ Each **circle** represents a **group  $j$**

- Pooled approach
- **Between approach**
- Random Effects
- **Fixed Effects**

>> Focus on within-group variation

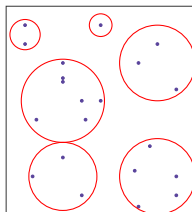
>> Implementation: dummy variables for each **group  $j$**  ( $\gamma_j$ )

$$Y_{ij} = \gamma_j + \beta \cdot X_{ij} + \varepsilon$$

>> What about panel data?

# Fixed-effect regression with panel data

>> Assume a pool of structured data



↪ Each **dot** represents a **measure**  $t$

↪ Each **circle** represents a **unit**  $i$

>> Focus on **within-unit** variation

>> Implementation: dummy variables for each **unit**  $i$  ( $\gamma_i$ )

$$Y_{it} = \gamma_i + \beta \cdot X_{it} + \varepsilon$$

>> What about time fixed-effect?



# DiD: Two-way fixed-effect regression

## Estimator (Regression with Multiple Time Periods)

*We can generalise to multiple groups/time periods using unit and period fixed-effects ('two-way' fixed-effect model):*

$$Y_{it} = \gamma_i + \alpha_t + \delta \cdot D_{it} + \varepsilon$$

- $\gamma_i$  is a fixed-effect for units (dummy for each unit)
- $\alpha_t$  is a fixed-effect for time periods (dummy for each period)
- $\delta$  is the DiD estimate based on  $D_{it}$

## Very flexible approach

- we can replace  $D_{it}$  with almost any type of treatment (not only binary)
- we can extend easily to multiple periods
- we can have units treated at different times
- we can estimate unit-specific time trends by including a unit-period interaction
  - ↪ useful when treatment occurs at different times for different units and there are slight deviations from parallel trends

# DiD: Two-way fixed-effect regression

## Why does two-way fixed-effect regression estimate the DiD?

~> Unit FEs means that we are only using *within unit* variation in  $Y$  to calculate the effect of  $D$

- i.e., *changes* over time!
- This removes *all* time-constant confounders

~> Time FEs means that we remove the effect of any changes to the response variable that affect all units at the same time

~>  $\hat{\delta} \rightarrow \hat{\delta}_{ATT}$  (it might not be that simple...)

- It is hard to provide a visual inspection of the parallel trends assumption here as treatment may switch on at different time for different units
- Nevertheless, we are still assuming that treated/control units would have evolved identically over time in absence of treatment

>> Why not always use unit dummies?

- Fine in panel data, as we have same units at several points in time
- Not possible with repeated cross-section when we do not have the same units in each time period

## Some caution with two-way fixed-effect models

---



**Maxim Ananyev**  
@maximananyev



A rare photo of an applied economist  
keeping up with the difference-in-  
differences literature



# Unit FEs and time-constant confounders



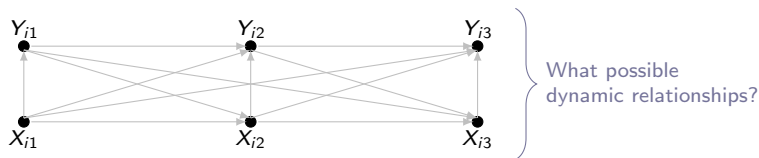
## When Should We Use Unit Fixed Effects Regression Models for Causal Inference with Longitudinal Data?

**Kosuke Imai**  
**In Song Kim**

Harvard University  
Massachusetts Institute of Technology

# Unit FEs and time-constant confounders

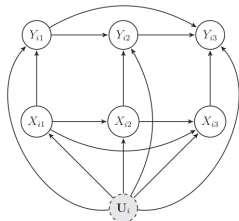
- ~ Imai & Kim (2019) show that unit FEs might not be that effective in adjusting for unobserved time-constant confounders
- ~ The issue is related to possible dynamic causal relationships



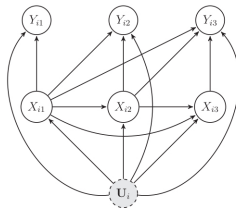
- ~ Some dynamic causal relationships compromise unit FEs

# Unit FEs and time-constant confounders

**FIGURE 2 Identification Assumptions of Regression Models with Unit Fixed Effects**



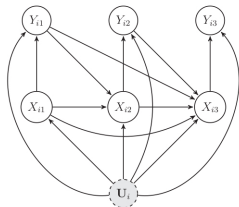
(a) past outcome affects current outcome



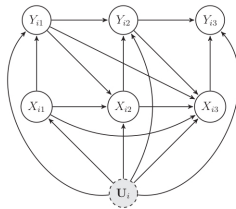
(b) past treatments affect current outcome

~ (1) Past outcome affects current outcome

~ (2) Past treatments affect current outcome



(c) past outcomes affect current treatment



(d) past outcomes affect both current outcome and treatment

~ (3) Past outcomes affect current treatment

~ (4) Past outcomes affect current outcome and treatment

# Unit FEs and time-constant confounders

## Key assumptions of unit fixed effects models

1. Past treatments do not directly influence current outcome
2. Past outcomes do not affect current treatment



# Summary

---

## Summary

- ~> Causal inference with observational data is really hard!
- ~> Longitudinal data can help, but it's not a silver bullet
  - Have a look at all assumptions involved
  - Parallel trends is an untestable assumption
- ~> This is a fast-changing topic. Keep up with the literature!
  - Callaway and Sant'Anna (2020); Callaway et al. (2021); Imai et al. (2021); Goodman-Bacon (2018); Imai and Kim (2019)
- ~> Now let's see how to estimate those models using R!
  - Find the lab notes here: [thiagoroliveira/2-LDA-lab.html](https://thiagoroliveira.com/2-LDA-lab.html)

# Thank you!

✉ [thiago.oliveira@manchester.ac.uk](mailto:thiago.oliveira@manchester.ac.uk)

🏠 [ThiagoROliveira.com](https://ThiagoROliveira.com)

🔗 [@oliveiratr.bsky.social](https://@oliveiratr.bsky.social)

## REFERENCES

- Callaway, B., A. Goodman-Bacon, and P. H. Sant'Anna (2021). Difference-in-differences with a continuous treatment. *arXiv preprint arXiv:2107.02637*.
- Callaway, B. and P. H. Sant'Anna (2020). Difference-in-differences with multiple time periods. *Journal of Econometrics*. Published online.
- Goodman-Bacon, A. (2018). Difference-in-differences with variation in treatment timing. Technical report, National Bureau of Economic Research.
- Imai, K. and I. S. Kim (2019). When should we use unit fixed effects regression models for causal inference with longitudinal data? *American Journal of Political Science* 63(2), 467–490.
- Imai, K., I. S. Kim, and E. H. Wang (2021). Matching methods for causal inference with time-series cross-sectional data. *American Journal of Political Science*.