

Tutorial 6: Introduction to Panel Data

Hammad Shaikh

March 4, 2021

Types of Data

④ Repeated Cross sectional: $(x_i, y_i)_{i=1}^I$ in $T \Rightarrow (x_j, y_j)_{j=1}^J$ in T'
↳ EX. Canadian census data

① • Cross sectional: same time period measured for different units
↳ $\{(x_i, y_i) : i=1, \dots, N\}$ in a single period
multiple units ↳ EX. Class grades for a given assignment

② • Time series: same unit measured at different time periods
↳ $\{(x_t, y_t) : t=1, 2, \dots, T\}$ for a single unit
↳ EX. Canadian GDP over years

③ • Panel data: range of units over time periods
↳ $\{(x_{it}, y_{it}) : i=1, 2, \dots, N, t=1, \dots, T\}$
NT data points ↳ EX. Daily Covid-19 cases for all Canadian Provinces

Panel Data Example

Table: Educational Attainment in Canada

Province	HS Graduation Rate	Years of Education	Year
Ontario	70	13	2000
⋮	⋮	⋮	⋮
Ontario	86.5	16	2018
⋮	⋮	⋮	⋮
Alberta	55	10	2000
⋮	⋮	⋮	⋮
Alberta	70	14	2018

- What are the variables? *HS Grad. and Years of Educ.*
- What is the time period? *Years from 2000-2018*
- What is the unit of observation? *Provinces*

Wide format

Prov	Educ2000	Educ2001	...
ON			
AL			
⋮			
⋮			

*Long format
because of
Year index*

** note: panel data
can be wide or
long format*

Regressions with panel data

Data: $\{(X_{it}, y_{it}) : i=1, \dots, N, t=1, \dots, T\}$

$\rightarrow \hat{\beta}_1$ unbiased if $\text{Cov}(X_{it}, \varepsilon_{it}) = 0$

- Pooled regression: $y_{it} = \beta_0 + \beta_1 x_{it} + \epsilon_{it}$

Ignoring the panel data structure

↳ Not leveraging either the time or unit index

- Individual fixed effects: $y_{it} = \beta_0 + \beta_1 x_{it} + \underbrace{\alpha_i + u_{it}}_{\varepsilon_{it}}$

↳ α_i represents any variables that affect y_{it} but are time-invariant and can differ across units

↳ Ex. Student's innate academic ability

Ex. Laws in a province in short time

First difference estimator

unit-level fixed effect

- Want to estimate: $y_{it} = \beta_0 + \beta_1 x_{it} + \alpha_i + u_{it}$
↳ Controls for any time-invariant unit level variables that affect outcome \Rightarrow (can handle endog. $\text{Cov}(x_{it}, \alpha_i) \neq 0$)
- First difference (FD) estimator: OLS of Δy_{it} on Δx_{it}

$$\underbrace{y_{it} - y_{it-1}}_{\Delta y_{it}} = \beta_1 \underbrace{(x_{it} - x_{it-1})}_{\Delta x_{it}} + \Delta u_{it} \Rightarrow \hat{\beta}_1 \text{ unbiased if } \text{Cov}(\Delta x_{it}, \Delta u_{it}) = 0$$

↳ no longer need to worry about $\text{Cov}(x_{it}, \alpha_i) \neq 0$

Difference - in - Differences (DD) for Policy Analysis

Video: <https://www.youtube.com/watch?v=V07MKhud-y0>

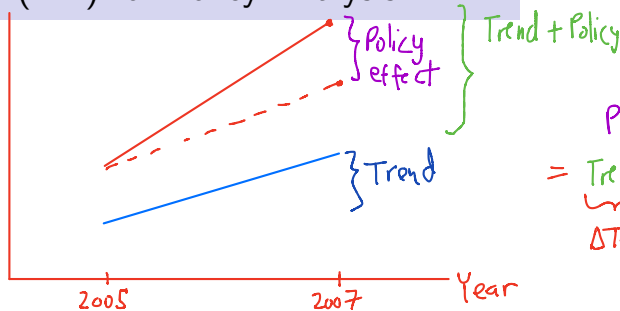
- DD used when have a pre-post and control-treatment setting (Environment)
Control units that never get policy
Treatment units that eventually get policy
Before & after policy
- Want to learn impact of funding on school performance (Policy)
↳ Policy is educational funding
- School performance in Toronto and Ottawa in 05 and 07 (Data)
Cities are units *Years are time periods*
- Suppose in 2007 Toronto received school funding (Treatment & Control)
↳ Toronto is the treatment city
Ottawa is control city

Difference - in - Differences (DD) for Policy Analysis

Ottawa

Toronto

School
Performance



- Visualize the analysis:

Assume that performance in
Ottawa is the same as Toronto
↳ "Common Trends"

- Turn visualization into a Diff.-in-Diff. regression:

$$Y_{ct} = \beta_0 + \beta_1 \text{Treat}_c + \beta_2 \text{PostTreat}_t + \beta_3 \text{Treat}_c \times \text{PostTreat}_t$$

$c = \text{city}$
 $t = \text{year}$

$\text{Treat}_c = \begin{cases} 1, \text{Toronto} \\ 0, \text{Ottawa} \end{cases}$

$\text{PostTreat}_t = \begin{cases} 1, \text{Year} = 2007 \\ 0, \text{Year} = 2005 \end{cases}$

$\Rightarrow \hat{\beta}_3$ is the policy effect
 ↳ p-in-D estimate

Panel data practice - Chapter 13, Question 5

$$\underbrace{S_{it}}_{\text{Savings}} = \beta_0 + \beta_1 \text{Year } 1992_t + \beta_2 \text{age}_{it} + \varepsilon_{it}$$

* Panel data
only include controls
that vary in
(i,t) index.

We want to estimate the effect of several variables on annual saving. Suppose we have a panel data set on individuals collected on 1990, and 1992. If we include a year dummy for 1992 and use first differencing, can we also include age in the original model? Explain.

$$\Delta S_{it} = S_{i92} - S_{i90} = \beta_1 (1) + \beta_2 \underbrace{\Delta \text{age}_{it}}_2 + \Delta \varepsilon_{it}$$

↳ Identification problem to separate β_1 from β_2

↳ Doesn't make sense to reg. ΔS_{it} on $\underbrace{\Delta \text{age}_{it}}_2$, no variation