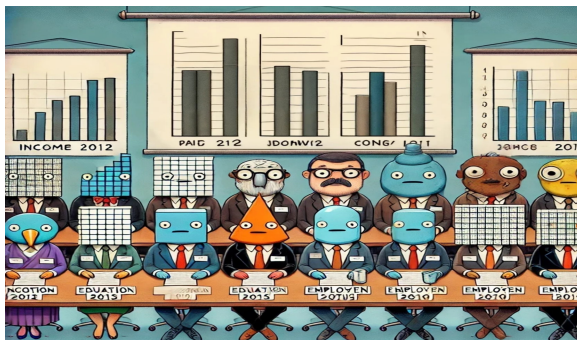


Fundamentals of Econometrics

Lecture 8: Pooling Cross Sections across Time: Simple Panel Data Methods



Section 1

Simple Panel Data Methods

- Panel data is a combination of time series and cross-sectional data.
- It consists of observations on multiple entities (individuals, firms, countries, etc.) over multiple time periods.
- To gather a panel data set, we can either:
 - Collect data on the same individuals over time (e.g., a survey of the same individuals at different points in time).
 - Collect data on different individuals at different points in time (e.g., a survey of different individuals at different points in time).

Independently Pooled Cross Sections

- Collection of independent, random samples from the same population at multiple periods of time.

Some Data Sources

- The Panel Study of Income Dynamics (PSID): Collected by the University of Michigan
- National Longitudinal Surveys (NLS): Collected by the Bureau of Labor Statistics
- Medical Expenditure Panel Survey (MEPS): Collected by the Agency for Healthcare Research and Quality
- National Health and Nutrition Examination Survey (NHANES): Collected by the National Center for Health Statistics at the CDC
- American Community Survey (ACS): Collected by the U.S. Census Bureau
- Current Population Survey (CPS): Collected by the Bureau of Labor Statistics
- American Time Use Survey (ATUS): Collected by the Bureau of Labor Statistics

Using Pooled Cross Sections Data

- We want to pool all cross-sections over time into one single data set

Advantages:

- 1 Increased sample size: By pooling data from multiple cross-sections, we can increase the sample size, which can lead to more precise estimates.
- 2 Improve generalizability: By pooling data from different time periods, we can improve the generalizability of our results to the population.
- 3 Can be used to estimate the effect of a policy change or event on a population (natural experiment).

Natural or Quasi-Experiments

- A natural experiment is when an exogenous shock occurs to a system and affects individual behavior.
- We have a group of individuals affected by the shock and a group of individuals that are not. So this is similar in principle to a laboratory experiment where there is a treatment group (affected by shock) and control group (not affected by shock).
- A quasi-experiment is when a researcher uses a natural experiment to estimate the effect of a treatment on an outcome.

Two requirements: - Two time periods (one before and one after the policy change) - Two groups (treatment and control)

Natural Experiment Framework

Goal: To determine differences between treatment and control groups due to an exogenous shock.

- 1 Pool the data from the two time periods.
- 2 Include a dummy variable for time and group.
 - $d2 = 1$ if obs occurs after event
 - $dT = 1$ if obs occurs in treatment group
- 3 Include additional variables and an interaction term between the two dummy variables

$$y = \beta_0 + \underbrace{\delta_0 d2}_{\substack{\text{Controls for} \\ \text{unobserved} \\ \text{changes affecting} \\ \text{both groups}}} + \underbrace{\delta_1 dT}_{\substack{\text{Controls for} \\ \text{initial difference} \\ \text{between groups}}} + \delta_2 d2 \cdot dT + \underbrace{\text{other factors}}_{\substack{\text{Controls for} \\ \text{observable differences} \\ \text{between treatment} \\ \text{and control group}}} + u$$

Group	Before	After	After - Before
Treatment	$\hat{\beta}_0 + \hat{\delta}_1$	$\hat{\beta}_0 + \hat{\delta}_0 + \hat{\delta}_1 + \hat{\delta}_2$	$\hat{\delta}_0 + \hat{\delta}_2$
Control	$\hat{\beta}_0$	$\hat{\beta}_0 + \hat{\delta}_0$	$\hat{\delta}_0$
{Treatment - Control}	$\hat{\delta}_1$	$\hat{\delta}_1 + \hat{\delta}_2$	$\hat{\delta}_2$

Example: Garbage Incinerator

- Blacksburg is considering the location of a new garbage incinerator. We want to conduct an analysis of the impact of the incinerator on property values.
- We could consider a town similar to Blacksburg where an incinerator was built. We will use Boston housing data from Kiel and McCain (1995). We have data from 1971 and in 1981 (when the incinerator was built).
- Similar town with:
 - 30,000 residents
 - Small college of 2,000 students
 - 25 miles from the nearest city

```

l.before <- lm(rprice ~ nearinc ,data = kielmc, subset = year == 1978)
l.after <- lm(rprice ~ nearinc ,data = kielmc, subset = year == 1981)

stargazer(l.before, l.after, font.size = "scriptsize",
  title = "Pooled Cross Sections", column.labels = c("Before", "After"),
  omit.stat = c("f", "ser", "adj.rsq"), header = FALSE)

```

Table 2: Pooled Cross Sections

	<i>Dependent variable:</i>	
	rprice	
	Before	After
	(1)	(2)
nearinc	-18,824.000*** (4,745.000)	-30,688.000*** (5,828.000)
Constant	82,517.000*** (2,654.000)	101,308.000*** (3,093.000)
Observations	179	142
R ²	0.082	0.165

Note: *p<0.1; **p<0.05; ***p<0.01

Garbage Incinerator Example

- What is the effect of the garbage incinerator on housing prices?
- Being close to an incinerator depresses prices, but location was one with lower prices to begin with.
- The effect of the garbage incinerator on housing prices is given by the difference in the coefficients of the `nearinc` variable in the two models.

$$\hat{\delta} = -30688.274 - (-18824.370) = \boxed{-11863.903}$$

$\hat{\delta}$ is known as the **difference-in-difference (DiD)** estimator. It can be expressed as the difference over time in the average difference in housing prices between the two groups (near and far from the incinerator):

$$\hat{\delta} = \left(\overline{rprice}_{81,near} - \overline{rprice}_{81,far} \right) - \left(\overline{rprice}_{78,near} - \overline{rprice}_{78,far} \right)$$

Final Words

Basic Setup of DiD:

- Two groups (treatment ($D_i = 1$) and control ($D_i = 0$))
- Two time periods (before ($T=0$) and after the treatment ($T=1$))

Group (D_i)	After Treatment ($T_i = 1$)	Before Treatment ($T_i = 0$)
Treated ($D_i = 1$)	$E[Y_{1i}(1) D_i = 1]$ $E[Y_{0i}(0) D_i = 1]$	
Control ($D_i = 0$) }	$E[Y_{0i}(1) D_i = 0]$	$E[Y_{0i}(0) D_i = 0]$

The **fundamental challenge**: We cannot observe $E[Y_{0i}(1)|D_i = 1]$ —i.e., the **counterfactual outcome** for the treated group had they not received treatment.

- DiD estimates the Average Treatment Effect on the Treated (ATT) as:

DiD Estimator in Regressions

$$rprice_{it} = \beta_0 + \delta_0 after + \beta_1 nearinc + \delta_1 after \cdot nearinc + u_{it}$$

- The differential effect of being in the location **and** after the incinerator was built is given by δ_1 .
- The DiD estimator is the difference in the coefficients of **nearinc** in the two models.
- If houses sold before and after the incinerator was built were systematically different, further explanatory variables should be included.
 - **This will also reduce the error variance and thus standard errors.**

DiD Estimator in Regressions

Table 4: Pooled Cross Sections

	<i>Dependent variable:</i>		
	rprice		
	(1)	(2)	(3)
y81	18,790.000*** (4,050.000)	21,321.000*** (3,444.000)	13,928.000*** (2,799.000)
nearinc	-18,824.000*** (4,875.000)	9,398.000* (4,812.000)	3,780.000 (4,453.000)
y81:nearinc	-11,864.000 (7,457.000)	-21,920.000*** (6,360.000)	-14,178.000*** (4,987.000)
Constant	82,517.000*** (2,727.000)	89,117.000*** (2,406.000)	13,808.000 (11,167.000)
Other Controls	No	Age, AgeSq	Full Set
Observations	321	321	321
R ²	0.174	0.414	0.660

Note:

*p<0.1; **p<0.05; ***p<0.01

```
l.difference <- lm(rprice ~ y81 + nearinc + y81:nearinc,  
                  data = kielmc)  
  
l.difference2 <- lm(rprice ~ y81 + nearinc + y81:nearinc +  
                   age + agesq,  
                   data = kielmc)  
  
l.difference3 <- lm(rprice ~ y81 + nearinc + y81:nearinc +  
                   age + agesq + intst + land + area +  
                   rooms + baths,  
                   data = kielmc)  
  
stargazer(l.difference, l.difference2, l.difference3,  
          font.size = "scriptsize",  
          title = "Pooled Cross Sections",  
          # Keep only y81, nearinc, y81:nearinc, "Constant")  
          keep = c("Constant", "y81", "nearinc", "y81:nearinc"),  
          add.lines = list(c("Other Controls", "No", "Age, AgeSq", "Full Set")),  
          omit.stat = c("f", "ser", "adj.rsq"), header = FALSE)
```


Table 5: Pooled Cross Sections

	<i>Dependent variable:</i>		
	log(rprice)		
	(1)	(2)	(3)
y81	0.193*** (0.045)	0.220*** (0.037)	0.162*** (0.028)
nearinc	-0.340*** (0.055)	0.007 (0.052)	0.032 (0.047)
y81:nearinc	-0.063 (0.083)	-0.185*** (0.068)	-0.132** (0.052)
Constant	11.300*** (0.031)	11.400*** (0.026)	7.650*** (0.416)
Other Controls	No	Age, AgeSq	Full Set
Observations	321	321	321
R ²	0.246	0.509	0.733

Note:

*p<0.1; **p<0.05; ***p<0.01

Model in Logs

```
1.difference.log <- lm(log(rprice) ~ y81 + nearinc + y81:nearinc,
  data = kielmc)
1.difference2.log <- lm(log(rprice) ~ y81 + nearinc + y81:nearinc +
  age + agesq,
  data = kielmc)
1.difference3.log <- lm(log(rprice) ~ y81 + nearinc + y81:nearinc +
  age + agesq + lintst + lland + larea +
  rooms + baths,
  data = kielmc)
stargazer(1.difference.log, 1.difference2.log, 1.difference3.log,
  font.size = "scriptsize",
  title = "Pooled Cross Sections",
  # Keep only y81, nearinc, y81:nearinc
  keep = c("Constant", "y81", "nearinc", "y81:nearinc"),
  add.lines = list(c("Other Controls", "No", "Age, AgeSq", "Full Set")),
  omit.stat = c("f", "ser", "adj.rsq"), header = FALSE)
```

Adding Additional Control Groups

- An implicit assumption of the DiD estimator is that the treatment and control groups are similar in all respects except for the treatment.
- We assume that potential trends in the outcome, y , would trend at the same rate in the absence of the treatment.
- This is known as the **parallel trends assumption**.
- If this assumption is violated, the DiD estimator will be biased.
- We can add an additional control group (say a state that has not received the treatment but has a similar trend pre-treatment) to the model.
- Leads us to a **triple difference** estimator (**Difference-in-Difference-in-Difference**).

Two Period Fixed Effects Model

- Panel data contains observations on the same individuals in every time period.
- The main advantage is that we can control for all of the unobserved features of individuals that do not change over time. We call these individual-specific characteristics either unobserved or fixed effects.

Two Period Fixed Effects Model:

$$y_{it} = \beta_0 + \gamma_0 d2 + \beta_1 x_{i1t} + \dots + \beta_k x_{ikt} + \theta_i + \varepsilon_{it}, \quad i = 1, \dots, N, t = 1, 2$$

where:

- $d2 = 1$ if the observation occurs after the event ($t=2$), and zero otherwise.
- θ_i is the unobserved fixed effect for individual i .

Two Period Fixed Effects Model

The composite error term is given by:

$$u_{it} = \theta_i + \varepsilon_{it}$$

In general, we need $E[u_{it}|x_{it}, \theta_i] = 0$.

But, OLS is unbiased if $E[\varepsilon_{it}|x_{it}, \theta_i] = 0$. - $Corr[X, \theta_i] \neq 0$ is okay - Can relax usual conditional mean assumptions for error term. - All comes out of measurement on the same individuals over multiple time periods.

First Difference Model

$$\text{Time 1: } y_{i1} = \beta_0 + \beta_1 X_{i11} + \dots + \beta_k X_{ik1} + \theta_i + \varepsilon_{i1}$$

$$\text{Time 2: } y_{i2} = (\beta_0 + \gamma_0) + \beta_1 X_{i12} + \dots + \beta_k X_{ik2} + \theta_i + \varepsilon_{i2}$$

Subtracting the first equation from the second:

$$y_{i2} - y_{i1} = \gamma_0 + \beta_1(X_{i12} - X_{i11}) + \dots + \beta_k(X_{ik2} - X_{ik1}) + (\varepsilon_{i2} - \varepsilon_{i1})$$

The FD Estimator is then:

$$\Delta y_i = \gamma_0 + \beta_1 \Delta X_{i1} + \dots + \beta_k \Delta X_{ik} + \Delta u_i$$

We will then: - Use data to calculate the first difference. - Estimate the model using OLS.

Disadvantages of the First Difference Model

- All time invariant variables are dropped from the model (e.g. binary variables like **gender**, **race**, etc.)
 - This is a problem if we want to include these variables in the model and to determine their effect on the dependent variable.
- If there is little variation in the independent over time the variance of the first difference will be small and the standard errors will be large.
 - This is a problem if we want to include these variables in the model and to determine their effect on the dependent variable.

Example: Traffic Fatalities

Do open container law decrease auto fatalities?

Dataset: `traffic1`

- 50 states plus DC
- `dthrte` = deaths per 100 million miles driven
- `open` = 1 if state has open container law, 0 otherwise
- `admn` = 1 if state can suspend license before DD conviction, 0 otherwise

Table 6: States with Open Container Law

	Deaths	Admn	Open
1985	137.6999999332428	21	19
1990	109.8999999141693	29	22


```
traffic1 |> summarize(  
  across(c(dthrte85, dthrte90, admn85, admn90, open85, open90),  
    sum)) |>  
  matrix(nrow = 2, dimnames = list(c("1985", "1990"),  
                                     c("Deaths", "Admn",  
                                       "Open")))) |>  
  knitr::kable(digits = 3, caption = "States with Open Container Law")
```

Model: $dthrte_{it} = \beta_0 + \gamma_t + \beta_1 open_{it} + \beta_2 admn_{it} + \theta_i + \varepsilon_{it}$

DiD Estimator:

$$dthrte_{it} = \beta_0 + \gamma_1 d2 + \beta_1 open_{it} + \delta_1 d2 \cdot open_{it} + \beta_2 admn_{it} + \delta_2 d2 \cdot admn_{it} + \varepsilon_{it}$$

```
lm(dthrte ~ d2*(open + admn), data = traff) |>  
summary()
```

```
##  
## Call:  
## lm(formula = dthrte ~ d2 * (open + admn), data = traff)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -0.795 -0.356 -0.115  0.330  1.715   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)    2.685      0.118   22.67  <2e-16 ***  
## d2             -0.590      0.176   -3.35   0.0011 **   
## open           -0.213      0.162   -1.31   0.1926   
## admn            0.230      0.160    1.44   0.1532   
## d2:open         0.272      0.229    1.19   0.2378   
## d2:admn        -0.169      0.227   -0.74   0.4590   
## ---
```

Key Takeaways

- We would conclude that open container laws have no effect on fatalities.
- Issues?
 - Causality: it may be that states enact open container laws because they have higher death rates.
 - So death rates may be causing a change in laws, not the other way around
 - Neither model accounts for much variation in fatalities. There are many unobserved variables that might be correlated with the open container laws. For example,
 - states with open container laws may also have stricter DUI laws, which may also reduce fatalities.
 - \$ amount of fines for speeding or other traffic infractions,
 - maintenance and safety of roads,
 - use of seat belts, etc.

```
traff <- traffic1 |> select(state:speed85) |>
  pivot_longer(cols = -state,
    names_to = c("var", "year"),
# chars before last 2 digits as var, last 2 as year
    names_pattern = "(.*)\\d{2})",
    values_to = "value") |>
  mutate(d2 = ifelse(year == "90", 1, 0)) |>
  pivot_wider(
    names_from = var,
    values_from = value
  )
```

FD Estimator: $\Delta dthrte_{it} = \beta_0 + \beta_1 \Delta open_{it} + \beta_2 \Delta admn_{it} + \varepsilon_{it}$

```
lm(dthrte90 - dthrte85 ~ I(open90 - open85) + I(admn90 - admn85), data = traffic1)
```

```
##  
## Call:  
## lm(formula = dthrte90 - dthrte85 ~ I(open90 - open85) + I(admn90 -  
##      admn85), data = traffic1)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -1.2526 -0.1434 -0.0032  0.1968  0.7968   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)    -0.4968    0.0524   -9.48  1.4e-12 ***  
## I(open90 - open85) -0.4197    0.2056   -2.04    0.047 *   
## I(admn90 - admn85) -0.1506    0.1168   -1.29    0.204   
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.344 on 48 degrees of freedom  
## Multiple R-squared:  0.119, Adjusted R-squared:  0.0819   
## F-statistic: 3.23 on 2 and 48 DF, p-value: 0.0482
```

What was the impact of changes in the open container law on traffic fatalities?

- Reduction in fatality for states that did not change the law?
- Reduction in fatality for states that changed the law?
- Reliability of the results?

Takeaways: Impacts of open contain law

- In 1985, there were an average of 2.7 deaths per 100 million miles driven
- For states that did not change their open container law:
 - 18.4% reduction in fatalities between 1985 and 1990 $(2.7 - 0.497)/2.7$
- For states that enacted open container law: 15.5% additional reduction in fatalities between 1985 and 1990 $(2.7 - 0.42)/2.7$

- Estimate true model
- Random sample from each cross-section
- Variation in each independent variable and no perfect collinearity
- $E[\varepsilon_{it}|x_{i1}, \dots, x_{iK}, a_i] = 0$ and $E[\Delta\varepsilon_{it}|\Delta x_{i1}, \dots, \Delta x_{iK}, a_i] = 0$ for all t
- Homoskedasticity: $var(\Delta\varepsilon_{it}|\Delta x_{i1}, \dots, \Delta x_{iK}) = \sigma^2 \quad \forall t$
- No autocorrelation: $cov(\Delta\varepsilon_{it}, \varepsilon_{is}|\Delta x_{i1}, \dots, \Delta x_{iK}) = 0$ for all $t \neq s$
- Normality of differenced errors

Unbiased with 1-4