Statistics II

Week 9:

Panel data and fixed effects

Lecture Review

Panel Data

As we introduced last week, a powerful way to achieve causal identification is to gather multiple observations over time.

Cross-sections: Samples of different units measured at the same point in time.

Time series: The same subject or unit measured at different points in time.

Panel data: Multiple units over multiple points in time.

Balanced panels, a.k.a. longitudinal data, record *the same* individuals over time in waves.

Two components of the panel data setup

(Remember slides from last week)

Cross section

Measure of different units, at one point in time.

CS	SU vote s	hares	
Unit	Y_{2014}	Y_{2020}	D
County A	42.1	38.5	0
County A County B	41.2	40.2	1

Static group comparison:

We compare the T and C groups after treatment.

$$Y_{i1}^{0} = \underbrace{\theta_{i}^{0}} + \underbrace{\delta_{1}} + \underbrace{\upsilon_{i1}^{0}}$$

$$Y_{i1}^{1} = \tau + \underbrace{\theta_{i}^{1}} + \underbrace{\delta_{1}} + \underbrace{\upsilon_{i1}^{1}}$$

Because observations are simultaneous, time period effect is the same and cancels out.

To estimate *tau*, we need to make strong **assumptions**:

Exogeneity: The idiosyncratic error is independent of treatment.

Random effects: unobserved differences in units are independent of treatment.

(Potential outcomes of control group are the same as the counterfactual potential outcomes for those being treated.)

That is assuming a lot...

Two components of the panel data setup

Temporal

Measure of one unit, at different points in time.

CSU vote shares				
Unit	Y_{2014}	Y_{2020}	D	
County A	42.1	38.5	0	
County B	41.2	40.2	1	

Longitudinal comparison:

We compare Y for the one unit before and after treatment.

$$Y_{i0}^{0} = \underbrace{\delta_{0}^{0}}_{0} + \underbrace{\theta_{i}}_{0} + \underbrace{\psi_{i0}^{0}}_{0}$$

$$Y_{i1}^{1} = \tau + \underbrace{\delta_{1}^{1}}_{0} + \underbrace{\theta_{i}}_{0} + \underbrace{\psi_{i1}^{1}}_{0}$$

Because we are comparing a unit to itself, unit-specific effects are the same across time and are ruled out.

To estimate tau, we need to make strong assumptions:

Exogeneity: The idiosyncratic error is independent of treatment. The effects of transitory forces cancel out over time.

Temporal stability: no impact of unobserved time-varying factors. In the absence of treatment, there would be no change in the mean of y.

(Average potential outcome for the unit observed does no change in time)

That is also assuming a lot...

Panel Data

Measure of **many** (and the same) units, at **many** different points in time.

CSU vote shares						
Unit	Y_{2014}	Y_{2020}	D			
County A	42.1	38.5	0			
County B	41.2	40.2	, 1			

With panel data we can relax certain assumptions. We **do not** need to assume

- Exogeneity: Because theta is constant with time-series data.
- Temporal stability: Because delta is constant in cross-sectional data.

We rely on a weaker assumption:

$$E[Y_{i1}^1 - Y_{i1}^0] - E[Y_{i0}^1 - Y_{i0}^0] = \tau + E[\epsilon_{i1}^1 - \epsilon_{i1}^0] - E[\epsilon_{i0}^1 - \epsilon_{i0}^0]$$

To identify *tau*, ____must be zero. For that to happen, the error terms of T and C *do not need to be the same* in each time period, they only need to change in parallel across time.

(Remember DiD parallel trends assumption)

Fixed Effects

$$Y_{it} = eta_0 + eta_1 D_{it} + heta_i + \delta_t + v_{it}$$

Treatment indicator ϵ_{it} The error term of the observation of one unit at one

point in time

In a panel data setup we can decompose the error term in:

$$\theta_i \quad \text{Captures unit fixed effects: unmeasured characteristics of the units that \textit{don't change in time} \text{ and do affect the outcome (y). Think of the size of a city, climate, location; also gender or racial identity.}$$

$$\delta_t \quad \text{Captures time fixed effects: effects that take place at a certain time period but affect the outcome (y) of all units simultaneously. Think of a global economic shock, changes in national government.}$$

$$\upsilon_{it} \quad \text{The 'idiosyncratic' error (classical error) that contains factors that are both specific to unit and time.}$$

With panel data we can cancel out both unit and fixed effects, even if we cannot observe or measure the variables involved.

We would only need to be careful for confounding variables that vary **both** by unit **and** time period.

Estimation (1/2)

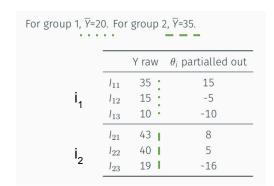
Fixed effects – or "de-meaning"

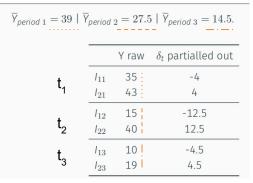
$\overline{ heta_i}$

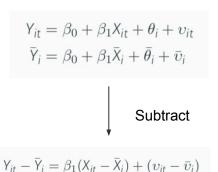
Given that the unit fixed effects are constant in time, we can remove their influence by subtracting the across time average outcome of each unit, from the y value of each observation.

δ_t

Likewise, given that time fixed effects are constant across units, we can remove their effect by subtracting the across unit average outcome of each period, from the y value of each observation.







If we estimate our regression model using this time-demeaned equation, we are left with the **FE estimator**: a model where all the confounders that don't vary over time just drop out.

Estimation (2/2)

Least Squares Dummy Variables (LSDV)

A second way to estimate fixed effects is to create dummy variables indicating the unit of every observation, and include them in the regression equation.

$$Y_{it} = \beta_0 + \beta_1 X_{it} + \underbrace{\beta_2 U_1 + \beta_3 U_2 + \dots + \beta_i U_{i-1}}_{unit \ dummies} + \upsilon_{it}$$

Two-way fixed effects: We can also add dummy variables for every time period observed to account for δ_t , unobserved time-specific effects. (This is a generalization of DiD).

In this case, do not include covariates that don't change over time in the model, or variables that only change over time but not across units.

In this case, do not include covariates that don't change over time in the model. They are already accounted for.

Questions?