

Causal Analysis using Panel Data

Tamer Çetin

Motivation

- You want to know how the industrial production of your country is affected by changes in the import demand of your country's largest trading partner.
- You have time series data on the industrial production of your country and total imports of its trading partner.
 - How should you estimate this effect?
 - Is there a way to get a reasonably precise effect estimate when your time series is not very long?
 - In particular, can you use similar time series from similar countries to get a good and more precise estimate of the effect for your country?
- You want to quantify the effects of immunization against measles.
- Among other potential effects, you want to know whether, and to what extent, immunization saves the lives of young children.
- You have access to data on immunization rates and child mortality from many countries from many years.
 - How should you use this data to get a good estimate of the effect?

Learning Outcomes

- After working through this chapter, you should be able to:
 - understand when and why to do pooled time series regressions to estimate effects, carry them out, and interpret their results;
 - carry out fixed-effects (FE) regression analysis on xt panel data and interpret its results;
 - carry out first-differenced (FD) regression analysis on xt panel data and interpret its results;
 - assess the extent to which the results of FE and FD regressions give good estimates of the average effect of a causal variable.

Time Periods

- Diff-in-diffs estimates the effect at a single point in time:
 - in the after, or endline, time period.
- When exactly we measure the effect is not important if the effect is immediate and stays the same.
- In most real-life situations, the effect may kick in after some delay, it may take some time to build up fully, and it may be more or less persistent, sometimes fading away in the longer run.
- In such cases, having a single endline time period is not enough to tell the full story.
- To be able to estimate how an effect plays out in time, we need more time periods.

Time Periods

- Another case when we need data from multiple time periods is when an intervention is scattered through time:
 - some subjects become treated at one point, some other subjects at another point, and so on.
- For example, when assessing the effect of the opening of large malls in small towns on the survival of small downtown shops, two-period diff-in-diffs can be applied if the mall openings were concentrated around the same time.
- In contrast, if malls opened in different years in different towns, scattered over a decade or more, the situation is unsuited for two-period diff-in-diffs.
- Instead, we would need data from small towns with many time periods, and we would need a more general method.

Estimating Effects Using Observational Time Series

- A time series regression can be specified in levels, with y_t regressed on x_t .
- However, we suggest to estimate time series regressions not in levels but in changes.
- A time series regression in changes, such as $\Delta y_t = \alpha + \beta \Delta x_t$, can uncover the effect of a change in x on how y changes within the same time period.
- A time series regression can uncover an average effect across the span of the times series for the single subject.
- For β to uncover the effect, Δx needs to be exogenous: time periods with different changes in x would have experienced the same change in y , had x changed the same way for them.

Lags to Estimate the Time Path of Effects

- One advantage of using data with multiple time periods is that we can estimate the time path of effects, such as immediate effects, effects in the near future, and **long-run effects**.
- To capture those in a **time series regression**, we need to include appropriate **lags** of Δx .
- With lags, we can estimate effects within the same time period (β_0 below), effects one time period later (β_1), and so on – provided, as usual, that x changes in an exogenous fashion.
- The time series regression that can estimate effects for up to K time periods has K lags of Δx :

$$\Delta y_t^E = \alpha + \beta_0 \Delta x_t + \beta_1 \Delta x_{t-1} + \dots + \beta_K \Delta x_{t-K}$$

Lags to Estimate the Time Path of Effects

- Then, by adding up the coefficients on all lags, we can estimate the long-run effect on y of a one-unit change in x .
- We can modify the right-hand-side variables to get a direct estimate of the long-run effect.

- With K lags of Δx , that is

$$\Delta y_t^E = \alpha + \beta_{cumul} \Delta x_{t-K} + \delta_0 \Delta(\Delta x_t) + \cdots + \delta_{K-1} \Delta(\Delta x_{t-(K-1)})$$

- where $\beta_{cumul} = \beta_0 + \beta_1 + \cdots + \beta_K$ above.
- β_{cumul} shows the **cumulative effect**:
 - the total change in y within K time periods after a unit change in x , on average.
 - When variation in Δx is exogenous, β_{cumul} shows the total effect of Δx_t on Δy over the long run.
 - We do not interpret the coefficients on the other variables here; in fact, we usually don't show them in the regression output.

Leads to Examine Pre-trends and Reverse Effects

- A time series regression, in differences, with K lags and L leads, has the form:
- $$\Delta y_t^E = \alpha + \beta_0 \Delta x_t + \beta_1 \Delta x_{(t-1)} + \cdots + \beta_K \Delta x_{(t-K)} + \gamma_1 \Delta x_{(t-1)} + \cdots + \gamma_L \Delta x_{(t-L)}$$
- The lag terms help capture delayed effects.
- The lead terms help capture differences in pre-trends and reverse effects.

Pooled Time Series to Estimate the Effect for One Unit

- When interested in the effect of x on y for a single unit, we can pool additional time series from similar units to obtain more precise estimates.

$$\Delta y_t^E = \alpha_i + \beta \Delta x_{it}$$

- It is good practice to allow for different intercepts for the different units while estimating a single slope coefficient to estimate an average effect.

Panel Regression with Fixed Effects

- The **pooled time series** is a useful starting point.
- One way to view panel data with multiple time periods is as if they were pooled time series.
- We can specify the regression in levels or changes, we can include leads and lags, and we should worry about trends, serial correlation, and seasonality.
- However, in typical multi-period panel data, we have many subjects (cross-sectional units) observed only a few times.
- This means that time series properties are less important, while differences across subjects are more important.
- Thus, it makes sense to discuss regression models in this context differently from the pooled time series setup.

Panel Regression with Fixed Effects

- The first regression model that we discuss for multi-period panel data is the **fixed-effects regression (FE regression)**.
- Fixed effects are separate intercepts in the regression for different cross-sectional units.
- In FE regressions we have y and x (in levels, not changes), and $xsec$ FE:
 - a separate intercept for each cross-sectional unit.
- The simplest linear panel regression with cross-sectional fixed effects is:

$$y_{it}^E = \alpha_i + \beta x_{it}$$

Panel Regression with Fixed Effects

- The fixed effects are denoted by α_i .
- These mean that the intercept is allowed to be different for different cross-sectional units.
- Why do we include the fixed effects?
- In other words, why do we include separate intercepts for each cross-sectional unit instead of including a common intercept?

Panel Regression with Fixed Effects

- The reason is the following.
- Suppose that subjects tend to have higher values of y on average due to some unobserved **confounder variable** that also affects x .
- Suppose, moreover, that the confounder affects x or y in the same way at all times.
- Then, with a common intercept α , we would have a usual linear regression as in cross-sectional data, and the coefficient β would be estimated with **omitted variables bias**.
- However, with the inclusion of the fixed effects, we can avoid, or mitigate, that bias.
- That's because including fixed effects means conditioning on all variables that don't change through time.

Panel Regression with Fixed Effects

- Technically, the inclusion of the cross-sectional **fixed effects** acts as a transformation of the y and x variables into differences from their cross-sectional means:
 - $y_{it} - \bar{y}_i$ and $x_{it} - \bar{x}_i$
where \bar{y}_i and \bar{x}_i are average values of y and x across all time periods within cross-sectional unit i .
- Indeed, it can be shown that the β in the model $y_{it}^E = \alpha_i + \beta x_{it}$ is exactly the same as the β in the model:
 - $(y_{it} - \bar{y}_i)^E - \alpha + \beta(x_{it} - \bar{x}_i)$.

Panel Regression with Fixed Effects

- As a result, in the FE regression, β shows how much larger y is, on average, compared to its mean within the cross-sectional unit, where and when x is higher by one unit compared to its mean within the cross-sectional unit.
- Or, to say it differently: compare two observations that are different in terms of the value of x compared to its i -specific mean.
- On average, y is larger, compared to its i -specific mean, by β , for the observation with the larger x value, compared to its i -specific mean.
- That is a **within-subject comparison**, and it is not affected by whether one subject has larger average y or x .
- That is why it is not affected by whether an unobserved confounder affects y or x in the same way in periods (and thus their average values).

Aggregate Trend

- The fixed-effects (FE) regression without covariates, using xt panel data with more than two time periods, is $y^E = \alpha_i + \theta_t + \beta x_{it}$
- Both x and y are in levels (as opposed to first differences).
- α_i are the cross-sectional FE; θ_t are the coefficients on time dummies, also known as time FE.
- β shows how much larger y is, on average, compared to its mean within the cross-sectional units and its mean within the time period, where and when x is higher by one unit compared to its mean within the cross-sectional unit and its mean within the time period.
- The cross-sectional FE regression can get us closer to estimating the effect of x on y by conditioning on confounders that do not change; with time dummies we can condition on aggregate trends of any shape.

Clustered Standard Errors

- When estimating regressions using xt panel data with more than two time periods, the standard errors need to be clustered at the level of the cross-sectional units.
- Clustered standard errors are robust to serial correlation as well as heteroskedasticity.

Panel Regression in First Differences

- Simple FD regression using xt panel data:

$$\Delta y^E = \alpha + \beta \Delta x_{it}$$

- β shows how much more y changes, on average, for observations with a one-unit higher increase in x .
- It is an average both across different cross-sectional units and different time periods.
- When Δx_{it} is exogenous, β shows the average effect of a change in x within the same time period.

Lags and Leads in FD Panel Regressions

- Panel regression in first differences with K lags:

$$\Delta y^E = \alpha + \beta_0 \Delta x_{it} + \beta_1 \Delta x_{i(t-1)} + \cdots + \beta_K \Delta x_{i(t-K)}$$

- The cumulative coefficient is

$$\beta_{cumul} = \beta_0 + \beta_1 + \cdots + \beta_K$$

- β_{cumul} shows how much more y tends to increase in total within the next K time periods, after a one-time increase of x by one more unit.
- When Δx is exogenous (contemporaneous and lagged), β_{cumul} shows the effect of a one-unit change of Δx on Δy , on average, over K time periods.
- We can estimate differences in pre-trends across cross-sectional units, and potential reverse effects, by including lead terms in the regression:

$$\Delta x_{i(t+1)}, \Delta x_{i(t+2)}, \cdots$$

Aggregate Trend and Individual Trends in FD Models

- FD panel regressions take care of aggregate linear trends automatically, by estimating associations of changes not levels.
- We can capture aggregate trends of any form by adding time dummies to FD regressions.
- We can capture individual linear trends by adding individual-specific intercepts to FD regressions.

Panel Regressions and Causality

- Suppose that the following regression would give a good estimate of the effect of x on y by conditioning on z :

$$y^E = \alpha + \beta x_{it} + \gamma z_i$$

- Importantly, here z_i differs across i but is the same for all t :
 - it does not change.
- That is in contrast to y_{it} and x_{it} , both of which differ not only across i but also across t .
- Suppose, moreover, that we do not observe z_i in the data.
- Recall our example of the effect of family income on fruit and vegetable consumption.

Panel Regressions and Causality

- Both may be affected by personality traits, such as how future-oriented people are.
- But those personality traits are unobserved in the data.
- Thus, regressing fruit and vegetable consumption on income confounds the effect of income and the effect of unobserved personality traits.
- Suppose, however, that personality traits don't change, at least not during the years we observe people in the data.
- That would be an example of a confounder that does not change.

Panel Regressions and Causality

- Now consider the FE regression.
- One way to understand the FE regression is that it transforms variables into mean differences:
 - their differences from the average value of each i across time.

- That's nothing other than the difference between the next two equations:

$$y^E = \alpha + \beta x_{it} + \gamma z_i$$
$$\bar{y}^E = \alpha + \bar{\beta} x_{it} + \bar{\gamma} z_i$$

- When we take the difference of these two equations, γz_i drops out because z_i does not change, and so its average value within i is simply its value:

$$z_i = \bar{z}_i$$

- So

$$(y_{it} - \bar{y}_i)^E = \beta (x_{it} - \bar{x}_i)$$

Panel Regressions and Causality

- But the same is true for FD regressions, which are differences between two time periods.
- Let us write out the original regression for times t and $t - 1$, and take their difference:

$$y^E = \alpha + \beta x_{it} + \gamma z_i$$

$$\Delta y_{i(t-1)}^E = \alpha + \beta x_{i(t-1)} + \gamma z_i$$

$$\Delta y_{it}^E = \beta \Delta x_{it}$$

First Differences or Fixed Effects?

- FD regressions:
 - can uncover immediate associations without including lags;
 - can uncover long-run associations by including sufficient numbers of lagged right-hand-side variables and calculating the cumulative slope coefficient;
 - take care of linear aggregate trends without including anything else;
 - can take care of nonlinear aggregate trends by including time dummies;
 - can take care of cross-sectional unit-specific linear trends by including cross-sectional FE.

First Differences or Fixed Effects?

- FE regressions:
 - can uncover long-run associations without including lags;
 - can uncover the time path of associations by including lags;
 - can take care of average nonlinear aggregate trends by including time dummies;
 - can take care of cross-sectional unit-specific linear trends in a cumbersome way.
- Summary advice:
 - FD is a better choice if interested in the time path of effects;
 - for long-run effect estimates, both FE and FD may be a good choice, with appropriate modifications.

Synthetic Control Method

- A comparative case study uncovers the effect of a single event, or intervention, on variable y for a single subject, by asking how the outcome would have changed without the event.
- The synthetic control method estimates the effect in a comparative case study framework.
- The synthetic control method creates a single control subject from a donor pool of many untreated subjects,
 - as a weighted average of the subjects in the donor pool;
 - by assigning weights to each subject in the donor pool;
 - making sure that the pre-intervention values of y and observed confounder variables are similar.

Synthetic Control Method

- The effect is estimated by comparing the observed y in the treated subject with the y value constructed from the synthetic control.
- Data analysts need to choose the donor pool, the pre-intervention variables, and the length of the pre-intervention time period.
- Then, the donors to the control group and their weights are selected by the algorithm.

Event Studies

- Event study regression

$$\Delta y_{it}^E = \alpha + \sum_0^{s_{max}} \beta_s D_{is} + \sum_{s_{min}}^1 \gamma_s D_{i(-s)}$$

- Event studies are a method to analyze the effect of an intervention (binary causal variable, or treatment) using xt panel data with subjects, some of which become treated during the time period covered in the data.
- Event studies re-define time around the intervention: this is called event time.
- Event studies provide a straightforward way to describe pre-intervention changes and post- intervention changes among treated subjects.

Selecting a Control Group in Event Studies

- Selecting a control group of untreated subjects is necessary to estimate the counterfactual:
 - What would have happened to treated subjects without the treatment.
- In event studies, we define the control group by defining pseudo-interventions.
- Pseudo-interventions are event time periods for untreated subjects that are preceded by changes in outcomes that are similar, on average, to pre-intervention changes in outcomes among treated subjects.

Main Takeaways

- Panel data methods help us get a step closer to causality.
 - Data with multiple time periods can help uncover short- and long-run effects and examine pre-trends.
 - When interested in the effects on a single cross-sectional unit, we may analyze a single time series or pool several time series of similar units.
 - With panel data having multiple time periods, we should use an FD regression to uncover the development of the effect over time, and an FD or an FE regression to uncover the long-run effect.

Main Takeaways

- When estimating the effect of an intervention using xt panel data, sometimes it's better to select a subset of non-treated observations to serve as a control group.
 - To estimate the effect of an intervention for a single subject, we can estimate the counterfactual using the synthetic control method.
 - With an intervention affecting many subjects at different times, we can carry out an event study with the help of a control group of comparable pseudo-interventions.

CASE STUDIES

- Import Demand and Industrial Production (Pooled Regressions)
- Immunization against Measles and Saving Children (FE and FD Regressions)
- Estimating the Effect of the 2010 Haiti Earthquake on GDP (Synthetic Control Method)
- Estimating the Impact of Replacing Football Team Managers (Event Study)