

Panel Data Models and Diff-In-Diff Methods

Zhiyuan Chen

RMEB

Remin Business School

April 2022

- Panel Data Models
 - Panel data
 - Random Effects vs. Fixed Effects
 - Fixed Effects vs. First Differencing
- Difference-In-Differences Models
 - Two-By-Two Case
 - Dynamic Treatment Effects and Event Study
 - Inference
- Staggered Difference-In-Differences
- Recent Developments in DID Models

- Now we have data for N units (people, products, firms, counties, cities, provinces, countries ...)
- For each unit we have T_i (≥ 1) observations
 - T_i typically refers to time periods
 - Could be products of a firm, students in a classroom, counties in a province

We assume that:

- N is large; consistency exists as N grows
- T_i is small; T_i stays constant as N grows

Now let's consider the regression model:

$$Y_{it} = X'_{it}\beta + u_{it}$$

- *Assumption 1*: assume independence **across individuals** but **not across time**
- *Assumption 2*: $\mathbf{E}(X_{it}u_{it}) = 0$

Consistency:

$$\hat{\beta} = \beta + \left(\frac{\sum_{i=1}^N \sum_{t=1}^{T_i} X_{it}X'_{it}}{\sum_{i=1}^N T_i} \right)^{-1} \left(\frac{\sum_{i=1}^N \sum_{t=1}^{T_i} X_{it}u_{it}}{\sum_{i=1}^N T_i} \right)$$
$$\rightarrow_p \beta$$

Serial Correlation

- The assumption of NO serial correlation is problematic so the asymptotic variance from before is not right
- Error terms are uncorrelated across individuals, but not within individuals
- Formally speaking, we assume that:

$$\mathbf{E}(u_{it}u_{jt}) = 0, \forall i \neq j$$

$$\mathbf{E}(u_{it}X_{it}) = 0$$

- However, the assumption of OLS is not satisfied because:

$$\mathbf{E}(u_{it}u_{is}) = 0$$

is a **CRAZY** assumption!

- Ignoring the serial correlation tends to **underestimate** the size of standard errors

- Recall that

$$\hat{\beta} - \beta = \left(\sum_i^N \sum_t^{T_i} X_{it} X'_{it} \right)^{-1} \sum_i^N \sum_t^{T_i} X_{it} u_{it}$$

- Define $\eta_i \equiv \sum_{t=1}^{T_i} X_{it} u_{it}$, then η_i is iid
- $\text{var}(\eta_i) = V_\eta$
- Then we get

$$\sqrt{N}(\hat{\beta} - \beta) = \left(\frac{1}{N} \sum_i^N \sum_t^{T_i} X_{it} X'_{it} \right)^{-1} \left(\frac{1}{\sqrt{N}} \sum_{i=1}^N \eta_i \right)$$

- Using the Central Limit Theorem we know

$$\sqrt{N}(\hat{\beta} - \beta) \sim \mathcal{N} \left(0, \left[\mathbf{E} \left(\sum_{t=1}^{T_i} X_{it} X'_{it} \right)^{-1} \right] V_\eta \left[\mathbf{E} \left(\sum_{t=1}^{T_i} X_{it} X'_{it} \right) \right]^{-1} \right)$$

We can approximate the variance of $\hat{\beta}$ as

$$\text{Var}(\hat{\beta}) \approx \left(\sum_i \sum_t X_{it} X'_{it} \right)^{-1} \left[\sum_i \left(\sum_t X_{it} \hat{u}_{it} \right) \left(\sum_t X'_{it} \hat{u}_{it} \right) \right] \left(\sum_i \sum_t X_{it} X'_{it} \right)^{-1}$$

- This is a generalization of the heteroskedastic robust standard errors
- We are allowing $[X_{i1}u_{i1}, X_{i2}u_{i2}, \dots, X_{iT_i}u_{iT_i}]$ to have an arbitrary variance-covariance matrix

In STATA, this is done by coding `reg y x, cluster(i)`

- i is the id of group in which the errors are serially correlated

Random Effects vs. Fixed Effects

- Panel data enable us to take care of the idiosyncratic component in the error term
- Write the linear model as

$$Y_{it} = X'_{it}\beta + \theta_i + \varepsilon_{it}$$

- In both of the models, we assume that

$$\mathbf{E}(\varepsilon_{it}X_{it}) = 0$$

- In the Random Effects (RE) model, we assume that

$$\mathbf{E}(\theta_i X_{it}) = 0$$

- In the Fixed Effects (FE) model, we do not assume anything about the relationship between θ_i and X_i ! θ_i is an idiosyncratic constant.

Consistent Estimators for FE models

- Include individual dummies and run the regression:

$$Y_{it} = X_{it}\beta + D'_{it}\theta + \varepsilon_{it}$$

where $D_{it} = [D_{it}^{(j)}]$, $D_{it}^{(j)} = 1$ if $i = j$ and zero otherwise

- Conceptually this is a different model, but technically it delivers consistent estimates for β
- **Standard FE (Mean-Differencing) model:**

$$(Y_{it} - \bar{Y}_i) = (X_{it} - \bar{X}_i)' \beta + (\varepsilon_{it} - \bar{\varepsilon}_i)$$

where $\bar{Z}_i \equiv \frac{1}{T_i} \sum_{t=1}^{T_i} Z_{it}$

- **First-Differencing (FD) model:**

$$Y_{it} - Y_{it-1} = (X_{it} - X_{it-1})' \beta + \varepsilon_{it} - \varepsilon_{it-1}$$

- With 2 periods, this is equivalent to the standard FE model (verify it by yourself)

Fixed Effects vs. First Differencing

More generally, we consider the model

$$Y_{it} = \beta \tau_{it} + \theta_i + \varepsilon_{it}$$

Assume that $T_i = T$ for everyone, for everyone the only regressor τ_{it} is given by

$$\tau_{it} = \begin{cases} 0 & t \leq t_0 \\ 1 & t > t_0 \end{cases}$$

- τ_{it} can represent some macro-level program (trade liberalization, 5-year plan, expansion of college enrollment...) begins at $t_0 + 1$
- The FE and FD estimators are:

$$FE : \hat{\beta}_{FE} = \frac{\sum_{i=1}^N \sum_{t=t_0+1}^T Y_{it}}{N(T-t_0)} - \frac{\sum_{i=1}^N \sum_{t=1}^{t_0} Y_{it}}{Nt_0}$$

$$FD : \hat{\beta}_{FD} = \frac{\sum_{i=1}^N (Y_{it_0+1} - Y_{it_0})}{N}$$

Fixed Effects vs. Pooled Regression

The fixed effects estimator is NOT always better than pooled OLS regression

- The sample variance of the data is:

$$\begin{aligned}ss(X_{it}) &= \sum_i \sum_t (X_{it} - \bar{X})^2 \\&= \underbrace{\sum_i \sum_t (X_{it} - \bar{X}_i)^2}_{\text{within variance}} + \underbrace{\sum_i T_i (\bar{X}_i - \bar{X})^2}_{\text{between variance}}\end{aligned}$$

FE estimator only uses the within variance of X_i :

- 1 It is inefficient. Standard errors are very large when X_{it} does not change
- 2 It can even a worse bias:

$$\begin{aligned}\hat{\beta}_{POLS} &= \beta + \frac{\text{cov}(X_{it}, \theta_i + \varepsilon_{it})}{\text{var}(X_{it})} \\ \hat{\beta}_{FE} &= \beta + \frac{\text{cov}(X_{it} - \bar{X}_i, \varepsilon_{it} - \bar{\varepsilon}_i)}{\text{var}(X_{it} - \bar{X}_i)}\end{aligned}$$

Example: Almond, Chay, and Lee (QJE, 2005)

- Their goal is to estimate the effects of birth weight on health:

$$h_{ij} = \alpha + \beta bw_{ij} + X_i' \gamma + a_i + \varepsilon_{ij}$$

where

- h_{ij} : health of newborn j of mother i
- bw_{ij} : birth weight
- a_i : mother specific effect

- The pooled OLS regression of h_{ij} on bw_{ij} gives us

$$\hat{\beta}_{POLS} = \beta + \frac{\text{cov}(bw_{ij}, X_i' \gamma)}{\text{var}(bw_{ij})} + \frac{\text{cov}(bw_{ij}, a_i)}{\text{var}(bw_{ij})}$$

- Their clever solution is to **use twins**:
 - Twins share the same mother, so a_i effectively controls for **race, age, education, family background...**
 - Estimate the model as

$$\Delta h_{ij} = \beta \Delta bw_{ij} + \Delta \varepsilon_{ij}$$

where they assume that $\text{cov}(\Delta bw_{ij}, \Delta \varepsilon_{ij}) = 0$

Difference-In-Differences

Difference-In-Differences: The Two-by-Two Case

Simple Before-After policy evaluation:

- Data on individuals right before and after the policy intervention: *Pre* and *Post*
- Two years dated 0 and 1 and that the policy is enacted in between
- We can simply identify the effect as:

$$\hat{\alpha} = \bar{Y}_1 - \bar{Y}_0$$

We could justify this using a FE model:

$$Y_{it} = \alpha_0 + \alpha\tau_{it} + \theta_i + \varepsilon_{it}$$

where

$$\tau_{it} = \begin{cases} 0 & t = 0 \\ 1 & t = 1 \end{cases}$$

We assume that $\mathbf{E}(\tau_{it}\varepsilon_{it}) = 0$ but no particularly assumption on θ_i

In the two-period case, FE model is equivalent to FD model, which is

$$Y_{i1} - Y_{i0} = \alpha + u_{i1} - u_{i0}$$

The estimator for the policy effect is

$$\hat{\alpha} = \bar{Y}_{i1} - \bar{Y}_{i0}$$

- We assume no other changes in time between and attribute whatever that is to the program
- Add time dummies into the model the treatment effect is not separated from the policy's impact

After adding time dummy variables, the model becomes

$$Y_{it} = \alpha_0 + \alpha \tau_{it} + \delta t + \theta_i + \varepsilon_{it}$$

Then the difference estimator delivers

$$\begin{aligned}\mathbf{E}(Y_{i1} - Y_{i0}) &= \mathbf{E}[(\alpha_0 + \alpha + \delta + \theta_i + \varepsilon_{i1}) - (\alpha_0 + \theta_i + \varepsilon_{i0})] \\ &= \alpha + \delta\end{aligned}$$

- The Diff-in-Diff estimator is designed to solve this problem

Now we have two periods for two groups:

- 1 People who are affected by the policy changes (treated): Y_{it}^1
- 2 People who are not affected by the policy changes (controls): Y_{it}^0

We can think:

- 1 Using the treated to pick up the time changes and policy effects:

$$\hat{\alpha} + \hat{\delta} = \bar{Y}_{i1}^1 - \bar{Y}_{i0}^1$$

- 2 Under the **common trend** assumption controls pick up the time changes:

$$\hat{\delta} = \bar{Y}_{i1}^0 - \bar{Y}_{i0}^0$$

We can then estimate the policy effect as a difference in difference

$$\hat{\alpha} = (\bar{Y}_1^1 - \bar{Y}_0^1) - (\bar{Y}_1^0 - \bar{Y}_0^0)$$

Formally, we can write the DGP (data generating process) for the DID estimator as

$$Y_{it} = \alpha_0 + \alpha \tau_{it}^{g_i} + \delta t + \theta_i + \varepsilon_{it}$$

where

$$\tau_{it}^g = \begin{cases} 1 & t = 1, g_i = 1 \\ 0 & \text{otherwise} \end{cases}$$

τ_{it}^g is usually written as $\tau_{it} \times t$

- Now we run a fixed effect regression and get a consistent estimate of the treatment effects α

Recall that the FE estimator for two-periods data is equivalent to the FD estimator:

$$\begin{aligned}\hat{\alpha} &= \frac{\sum_{i=1}^N \left[(\tau_{i1}^{g_i} - \tau_{i0}^{g_i}) - \overline{(\tau_{i1}^{g_i} - \tau_{i0}^{g_i})} \right] (Y_{i1} - Y_{i0})}{\sum_{i=1}^N \left[(\tau_{i1}^{g_i} - \tau_{i0}^{g_i}) - \overline{(\tau_{i1}^{g_i} - \tau_{i0}^{g_i})} \right]^2} \\ &= (\bar{Y}_1^1 - \bar{Y}_0^1) - (\bar{Y}_1^0 - \bar{Y}_0^0)\end{aligned}$$

- Notice that

$$\begin{aligned}\overline{(\tau_{i1}^{g_i} - \tau_{i0}^{g_i})} &= \frac{1}{N_1 + N_0} \sum_i (\tau_{i1}^{g_i} - \tau_{i0}^{g_i}) \\ &= \frac{1}{N_1 + N_0} \left[\sum_{\{i: g_i=1\}} (\tau_{i1}^1 - \tau_{i0}^1) + \sum_{\{i: g_i=0\}} (\tau_{i1}^0 - \tau_{i0}^0) \right] \\ &= \frac{N_1}{N_1 + N_0}\end{aligned}$$

In principle, you don't need panel data to implement DID regression; repeated cross section data are fine!

- In general, we can write the regression as

$$Y_i = \alpha_0 + \alpha \tau_{t_i}^{g_i} + \delta t_i + \gamma g_i + \varepsilon_i$$

- where g_i is the group indicator of person i and t_i is the time period in which i exists in the data
- $\tau_{t_i}^{g_i} = g_i \times t_i$
- We have 4 categories of people: (before, treated), (before, controls), (after, treated), (after, controls)
- We end up with having 4 groups of *moment conditions*:

$$\mathbf{E}(\varepsilon_i) = 0$$

$$\mathbf{E}(\tau_{t_i}^{g_i} \varepsilon_i) = 0$$

$$\mathbf{E}(t_i \varepsilon_i) = 0$$

$$\mathbf{E}(g_i \varepsilon_i) = 0$$

Using

$$Y_i = \hat{\alpha}_0 + \hat{\alpha}\tau_{t_i}^{g_i} + \hat{\delta}t_i + \hat{\gamma}g_i + \hat{\varepsilon}_i$$

We can rewrite the *sample analog* of the moment conditions as

$$\bar{Y}_0^1 = \hat{\alpha}_0 + \hat{\gamma}$$

$$\bar{Y}_1^1 = \hat{\alpha}_0 + \hat{\alpha} + \hat{\delta} + \hat{\gamma}$$

$$\bar{Y}_0^0 = \hat{\alpha}_0$$

$$Y_1^0 = \hat{\alpha}_0 + \hat{\delta}$$

We can solve for the parameters as

$$\hat{\alpha}_0 = \bar{Y}_0^0$$

$$\hat{\gamma} = \bar{Y}_0^1 - \bar{Y}_0^0$$

$$\hat{\delta} = \bar{Y}_1^1 - \bar{Y}_0^1$$

$$\hat{\alpha} = (Y_1^1 - \bar{Y}_0^1) - (\bar{Y}_1^0 - \bar{Y}_0^0)$$

More generally, we can add more control variables and specify the DID model as

$$Y_i = \alpha \tau_{t_i}^{g_i} + X_{it}' \beta + \delta_{t_i} + \theta_{g_i} + \varepsilon_i$$

This setting is simple yet powerful. There are many papers that do this basic stuff.

- Eissa and Liebman (QJE, 1996): Estimate the effect of the earned income tax credit on labor supply of women
- Dohahue and Levitt (QJE, 2001): Impact of legalized abortion on crime

- The treatment effects could be dynamic, with short-run effects differing from long-run effects
- We can easily extend the baseline model to allow for this:

$$Y_i = \beta_0 + \sum_{j=0}^J \alpha_j \tau_{t_i-j}^{g_i} + \delta_{g_i} + \gamma_{t_i} + \varepsilon_i$$

where

$$\tau_t^g = \begin{cases} 1 & g_i = 1 \text{ and policy started in year } t \\ 0 & \text{otherwise} \end{cases}$$

Common Trend Assumption

- Recall that the model we specified for DID is:

$$Y_i = \beta_0 + \alpha \tau_{t_i}^{g_i} + \delta t_i + \gamma g_i + \varepsilon_i$$

The DID estimator is

$$\hat{\alpha}_{DID} = \alpha + (\bar{\varepsilon}_1^1 - \bar{\varepsilon}_0^1) - (\bar{\varepsilon}_1^0 - \bar{\varepsilon}_0^0)$$

To obtain a consistent estimator, we require that

$$\mathbf{E} [(\bar{\varepsilon}_1^1 - \bar{\varepsilon}_0^1) - (\bar{\varepsilon}_1^0 - \bar{\varepsilon}_0^0)] = 0$$

- The treated units can have different *levels* of the error term, but the *change* in the error term must be random
- Two approaches to validate the common trend assumption: **placebo policies** and **add group-specify time trends**

- A popular strategy for robustness check: if a policy was enacted in say 2010 you could pretend it was enacted in 2005 in the same place and then only use data through 2009
- The easiest (and most common) is in the Event framework: include leads as well as lags in the model (Bertrand, Duflo, Mullainathan, 2004, QJE)
- To implement it you can run a regression like:

$$Y_i = \beta_0 + \sum_{j=-\tilde{J}, j \neq -1}^J \alpha_j \tau_{t_i-j}^{g_i} + \delta_{g_i} + \rho_{t_i} + \varepsilon_i$$

where $t_i = -\tilde{J}, \dots, J$ and the period of enactment to be zero

- $\tau_{t_i-j}^1 = 1$ whenever $t_i - j = 0$
- Note $\sum_{j=-\tilde{J}}^J \tau_{t_i-j}^{g_i} = \delta_{g_i}$, so we drop one period to avoid perfect co-linearity

If we normalize $\alpha_{-1} = 0$, α_{-2} is estimated as

$$\hat{\alpha}_{-2} = (\bar{Y}_{-2}^1 - \bar{Y}_{-1}^1) - (\bar{Y}_{-2}^0 - \bar{Y}_{-1}^0)$$

where

$$\bar{Y}_{-2}^1 = \hat{\beta}_0 + \hat{\alpha}_{-2} + \hat{\delta} + \hat{\rho}_{-2}$$

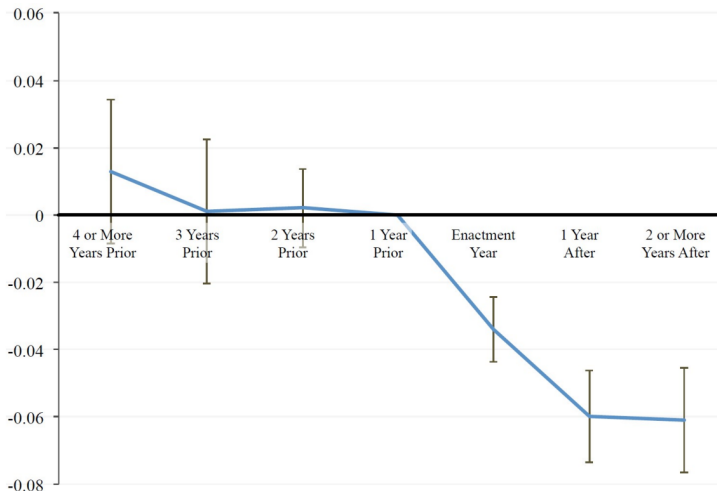
$$\bar{Y}_{-1}^1 = \hat{\beta}_0 + \hat{\delta} + \hat{\rho}_{-1}$$

$$\bar{Y}_{-2}^0 = \hat{\beta}_0 + \hat{\rho}_{-2}$$

$$\bar{Y}_{-1}^0 = \hat{\beta}_0 + \hat{\rho}_{-1}$$

- $\hat{\alpha}_{-2}$ should be zero under the common trend assumption
- $\hat{\alpha}_{-j} (j \geq 2)$ should also be zero under the common trend assumption

Figure 1: Typical event studies for the placebo policies



Group-Specific Time Trends

- One might be worried that units that are trending up or trending down are more likely to change policy
- One can include group \times time dummy variables in the model to fix this problem
- Now let's consider the three-period case where the policy happens between period 1 and 2, the model becomes

$$Y_i = \beta_0 + \alpha \tau_{t_i}^{g_i} + \delta_0 t_i \times (1 - g_i) + \delta_1 t_i \times g_i + \delta_2 \mathbb{I}(t_i = 2) + \gamma g_i + \varepsilon_i$$

- We can write the estimator of α as a Triple Difference (DDD):

$$\begin{aligned}\hat{\alpha}_{DDD} &= (\bar{Y}_2^1 - \bar{Y}_1^1) - (\bar{Y}_2^0 - \bar{Y}_1^0) - [(\bar{Y}_1^1 - \bar{Y}_0^1) - (\bar{Y}_1^0 - \bar{Y}_0^0)] \\ &\approx (\alpha + \delta_1 + \delta_2) - (\delta_0 + \delta_2) - (\delta_1 - \delta_0) \\ &= \alpha\end{aligned}$$

If we do not add the group-specific time trends, the error term is a composite:

$$\tilde{\varepsilon}_i = \varepsilon_i + \delta_0 t_i \times (1 - g_i) + \delta_1 t_i \times g_i - \delta t_i$$

We go back to the common trend assumption and get

$$\begin{aligned} \mathbf{E}((\tilde{\varepsilon}_1^1 - \tilde{\varepsilon}_1^0) - (\tilde{\varepsilon}_1^0 - \tilde{\varepsilon}_0^0)) &= \mathbf{E}((\bar{\varepsilon}_1^1 - \bar{\varepsilon}_1^0) - (\bar{\varepsilon}_1^0 - \bar{\varepsilon}_0^0)) \\ &\quad + (\delta_1 - \delta) - (\delta_0 - \delta) \\ &= \mathbf{E}((\bar{\varepsilon}_1^1 - \bar{\varepsilon}_1^0) - (\bar{\varepsilon}_1^0 - \bar{\varepsilon}_0^0)) + \delta_1 - \delta_0 \\ &\approx \delta_1 - \delta_0 \end{aligned}$$

Inference: Get the Right Standard Errors

- In applications, we often have individual data (people, firms...) and the policy is enacted at more aggregated level (cities, provinces)
- It seems that we have more data than we need

Consider the individual-level model

$$Y_i = \alpha \tau_{t_i}^{g_i} + Z_i' \delta + X_{g_i t_i}' \beta + \theta_{g_i} + \gamma_{t_i} + u_i$$

One could aggregate the data by group and year:

$$\bar{Y}_{gt} = \alpha \bar{\tau}_{gt} + Z_{gt}' \delta + X_{gt}' \beta + \theta_g + \gamma_t + \bar{u}_{gt}$$

The second regression gives consistent estimates for the treatment effects

- **HOWEVER, two regressions yield different standard errors: which one should we use?**

- The error term for the individual level model is

$$u_i = \eta_{giti} + \varepsilon_i$$

- The i.i.d assumption is violated is η_{giti} exists
- Standard solution: *Cluster the standard error by group \times year* to allow for arbitrary correlation within a group
- Bertrand, Duflo, and Mullainathan (QJE, 2004):
 - If there is serial correlation in η_{gt} , cluster by group \times year is problematic
- *Solutions:*
 - **# of groups is large:** block bootstrap by state works well
 - **Moderate # of states:** asymptotic approximation of the variance-covariance matrix
 - **Small # of states:** collapse the data into a “pre”- and “post”-period and account for the effective sample size
- A practioner’s guide is offered by Cameron and Miller (2015, JHR)

Extensions of DID

PSM-DID

Potential Outcome Framework

General Idea of Matching

Think of whether taking Chen's RMEB class as a binary decision, D_i ; and income in the future Y_i is the interested outcome:

$$\text{Potential Outcome} = \begin{cases} Y_i^1, & \text{if } D_i = 1 \\ Y_i^0, & \text{if } D_i = 0 \end{cases}$$

We aim to estimate

$$ATT = \mathbf{E}(\textcolor{blue}{Y}_i^1 - \textcolor{red}{Y}_i^0 | D_i = 1, \mathbf{X}_i), \quad D_i \in \{0, 1\}$$

The problem is only $Y_i = D_i Y_i^1 + (1 - D_i) Y_i^0$ is observed. In lack of ideal experiment, we use nonparticipants (control group) to approximate participants:

$$\underbrace{E(Y_i | D_i = 1, \mathbf{X}_i) - E(Y_i | D_i = 0, \mathbf{X}_i)}_{\text{observed income difference}} = \underbrace{E(Y_i^1 - Y_i^0 | D_i = 1, \mathbf{X}_i)}_{\delta(\mathbf{X}_i) = ATT | \mathbf{X}_i} + \underbrace{E(Y_i^0 | D_i = 1, \mathbf{X}_i) - E(Y_i^0 | D_i = 0, \mathbf{X}_i)}_{\text{Selection Bias}}$$

Conditional Independence Assumption

General Idea of Matching

- **Selection Bias:**

$$Bias(\mathbf{X}_i) = \mathbf{E}(Y_i^0 | D_i = 1, \mathbf{X}_i) - \mathbf{E}(Y_i^0 | D_i = 0, \mathbf{X}_i)$$

- **Conditional Independence Assumption (CIA):** $\{Y_i^0, Y_i^1\} \perp D_i | \mathbf{X}_i$
 - Under CIA, $Bias(\mathbf{X}_i) = 0$
 - CIA essentially states that D_i is randomly assigned conditioning on observable characteristics \mathbf{X}_i (see Angrist (1998) on voluntary military service)
 - CIA is likely to fail when *unobserved characteristics* determine select-into-treatment

“The idea of matching between treated and untreated units assumes that $Bias(\mathbf{X}_i) = 0$ so that conditioning on \mathbf{X}_i eliminates the bias.”

Regression vs. Matching

General Idea of Matching

- Statisticians more often use matching methods (Cochrane and Rubin, 1973); the idea is quite similar to conditioning on observables in the regression analysis
- “Regression can be motivated as a particular sort of weighted matching estimator” (Angrist and Pischke: *Mostly Harmless Econometrics*).
 - ① **Matching estimator** puts the most weight on covariate cells containing units who are most likely to be treated: high $\Pr(D_i = 1|\mathbf{X}_i)$
 - ② **Regression** puts the most weight on covariate cells with the largest conditional variance of treatment status: when $\Pr(D_i = 1|\mathbf{X}_i) \times [1 - \Pr(D_i = 1|\mathbf{X}_i)]$ is large
 - ③ **Common Support**: No weights assigned to covariate cells containing no treated and control units: $0 < \Pr(D_i = 1|\mathbf{X}_i) < 1$

Propensity Score Matching

A Structural Interpretation by Heckman, Ichimura, and Todd (1998)

- The dimension of \mathbf{X}_i can be high, adding to difficulty of matching
- The well-known PSM is to match over the propensity of selection into the treatment
- **Exclusion Restrictions:** Partition \mathbf{X} into (\mathbf{T}, \mathbf{Z}) such that:

$$Y_i^0 = g_0(\mathbf{T}_i) + U_i^0, \quad (1)$$

$$Y_i^1 = g_1(\mathbf{T}_i) + U_i^1, \quad (2)$$

$$\Pr(D_i = 1 | \mathbf{X}_i) = \Pr(D_i = 1 | \mathbf{Z}_i) = P(\mathbf{Z}_i) \quad (\text{propensity score}) \quad (3)$$

\mathbf{T} and \mathbf{Z} are not necessarily mutually exclusive

- To identify the ATT, it's enough to assume

$$U_i^0 \perp D_i | P(\mathbf{Z}_i) \quad (\text{PSM CIA})$$

Propensity Score Matching

A Structural Framework by Heckman, Ichimura, and Todd, (1998)

PSM method can be put in the context of classical econometric selection models:

$$Y_i^0 = g_0(\mathbf{Z}_i) + U_i^0 \quad (4)$$

$$D_i = \begin{cases} 1 & \text{if } \lambda(\mathbf{Z}_i) - v \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

- If \mathbf{Z}_i and v are independent, then $P(\mathbf{Z}_i) = \Pr\{v \leq \lambda(\mathbf{Z}_i)\} = F_v[\lambda(\mathbf{Z}_i)]$; the CIA implies that

$$E\{U_i^0 | D_i = 1, F_v[\lambda(\mathbf{Z}_i)]\} = E\{U_i^0 | D_i = 0, F_v[\lambda(\mathbf{Z}_i)]\}$$

- If also $\lambda(\mathbf{Z}_i) \perp (U_i^0, v)$ and $E(U_i^0) = 0$, then for any s :

$$E(U_i^0 | v = s) = 0 \quad (\text{No Selection on Observables})$$

Semi-parametric DID

Propensity Score Matching based DID Estimator (Heckman et al. (1998, ReSTud), Abadie (2005, ReSTud))

- Recall the OLS estimator the two-by-two DID model

$$Y_{it} = \delta D_{it} + \lambda_t + \eta_i + v_{it}$$

is simply

$$\delta = E[Y_{i1} - Y_{i0} | D_{i1} = 1] - E[Y_{i1} - Y_{i0} | D_{i1} = 0]$$

A sufficient condition for the identification is

$$\Pr(D_{i1} = 1 | v_{it}) = \Pr(D_{i1} = 1)$$

- If v_{it} is correlated in time, the model cannot be identified: What are possible solutions?

- The traditional solution is to add covariates \mathbf{X}_i into the model:

$$Y_{it} = \mu + \tau D_{it} + \mathbf{X}_i' \alpha(t) \lambda_t + \varepsilon_{it}$$

- Synthetic control estimator by Abadie, Diamond, and Hainmueller (2010, JASA):

$$\hat{\alpha}_t = Y_t^1 - \sum_{i=2}^{I+1} \omega_i^* Y_{it}^*$$

- $\omega^* = (\omega_2^*, \dots, \omega_{I+1}^*)$ is chosen to minimize $\|X_1 - X_0 \omega\|$
- Synthetic Difference in Differences estimator by Arkhangelsky et al. (2019, NBER working paper):

Nonlinear DID (Change-In-Changes)

Staggered Diff-In-Diffs

Appendices