

# Difference-in-Differences

Alessandro Caggia

June 2025

## Abstract

### 1. PANEL DATA

#### 2. 1.1. Idea

3 Now we have **panel data**, so we can have repeated measures for a  
4 given individual.

- **Issue:** If **unobserved variables** (selection on observables not possible) are correlated with treatment ( $D$ ) and potential outcomes ( $Y_{0i}, Y_{1i}$ )<sup>1</sup>, **unconfoundedness will not hold**:

$$(Y_{0i}, Y_{1i}) \not\perp D_i$$

5 e.g., **smarter people** (unobserved IQ, higher ( $Y_{0i}, Y_{1i}$ )) **self-select**  
6 into a supplementary math course: T will have a higher score  
7 than C in the math exam, but this is not driven entirely by the  
8 supplementary course but by intrinsic ability.

- **Solution:** The goal of Difference-in-Differences (DID) and panel data models is to leverage **panel data** to control for **permanent unobserved factors** (e.g., individual fixed effects): if we have more than one observation per individual, we can account for **time-invariant unobserved variables**.

- In a cross section: if someone has a high outcome, is it due to smartness or the treatment? We cannot tell as there are no past observations. Since smartness is unobserved, we cannot control it directly.
- With repeated observations, we can distinguish the **endogenous time-invariant smartness component** (pre-existing higher outcomes) from the treatment effect, and use the between-group comparison to disentangle the treatment effect from time effects.

- Go from the individual to the group level: imagine you have two groups. First, you remove the **endogenous time-invariant component** in the two groups exploiting **within-group variation**. Second, you exploit **between-group variation** to estimate the TE.

27 <sup>1</sup> The problem is not heterogeneous TEs, the problem is selection into treatment. If individuals have heterogeneous TEs, and this heterogeneity is equally represented across treated and control groups, it does not bias the estimate in fact, it is exactly the TE we want to identify.

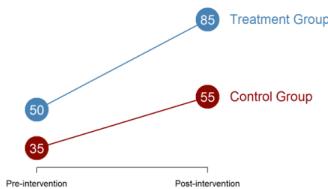


Figure 1. 2 time periods

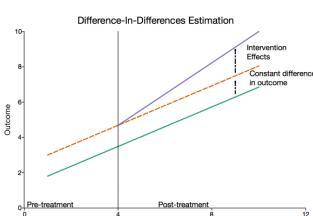


Figure 2. n time periods

#### 31 1.2. Unconfoundedness

We now assume unconfoundedness holds if we control for a time trend, observed time-varying covariates, and a vector of time-invariant unobserved variables  $A_i$  (expansions with respect to our

standard model):

$$E(Y_{0it} | A_i, X_{it}, t, D_{it}) = E(Y_{0it} | A_i, X_{it}, t)$$

#### 1.3. Functional form

- We impose a linear, additive functional form for potential outcomes:

$$E(Y_{0it} | A_i, X_{it}, t) = \alpha + \lambda_t + A_i'\gamma + X_{it}'\beta$$

$$E(Y_{1it} | A_i, X_{it}, t) = E(Y_{0it} | A_i, X_{it}, t) + \rho D_{it}$$

- Assume that Treatment Effects are homogeneous and constant over time:  $\rho_{it} = \rho$ .

- The switching equation:

$$\begin{aligned} E(Y_{it} | A_i, X_{it}, t, D_{it}) &= E(Y_{0it} | A_i, X_{it}, t) + D_{it}E(Y_{1it} - Y_{0it} | A_i, X_{it}, t) \\ &\quad \text{CME: D out} \\ &= \alpha + \lambda_t + A_i'\gamma + X_{it}'\beta + \rho D_{it} \end{aligned}$$

- The regression equation at the individual level is then (above is Expected Val):

$$Y_{it} = \alpha_i + \lambda_t + X_{it}'\beta + \rho D_{it} + \varepsilon_{it}$$

- where  $\alpha_i \equiv \alpha + A_i'\gamma$  is the individual fixed effect (time-invariant heterogeneity (e.g., ability, preferences)). If not controlled for, it creates bias:

##### – Composition:

- \*  $\alpha$ : common intercept across individuals.

- \*  $A_i'\gamma$ : a vector of time-invariant unobserved variables.

##### – Identification: $\alpha$ and $A_i'\gamma$ not separately identified.

- \* One intercept absorbs both  $\alpha$  and  $A_i'\gamma$ .

- \* Need to impose restrictions (eg, if  $A_i$  is normally distributed one could say that alpha is the mean of  $\alpha_i$ )

- Autocorrelation: this structure induces serial correlation in  $u_{it} = \alpha_i + \varepsilon_{it}$  for a given individual, since  $\alpha_i$  is shared across all  $t$ . The issue is generated by the  $A_i'\gamma$  component, not by the general constant, as the general constant shifts equally all individuals.

- $\varepsilon_{it}$  is the **idiosyncratic error**, varying over time and individuals.

#### 1.4. How to get rid of $\alpha_i$

##### 1.4.1. Random Effects

- **Random Effects:**  $\alpha_i$  (unobserved) is moved into the error term, assuming it is uncorrelated with the regressors.

- We can use pooled OLS, pooling all time periods. Better: for efficiency, use GLS, which accounts for correlation in errors (again,  $u_{it} = \alpha_i + \varepsilon_{it}$ ).

- If errors are correlated with regressors, you will have bias.
- REs are no longer used; they were more common in cases where you have low observations or low computing power.

##### 1.4.2. Fixed Effects and First Differences

- **Fixed Effects:**  $\alpha_i$  and  $\lambda_t$  are treated as parameters to estimate (via individual and time dummies). FE allows  $\alpha_i$  to correlate with regressors and controls for it;

- Using individual dummies generates a demeaned model (**you are controlling at the individual level, exploiting within-individual (over-time) variation, by removing individual-specific means**):

$$Y_{it} - \bar{Y}_i = \lambda_t - \bar{\lambda} + (X'_{it} - \bar{X}'_i)\beta + \rho(D_{it} - \bar{D}_i) + \varepsilon_{it} - \bar{\varepsilon}_i$$

- The rationale is: you remove the average, so necessarily you remove the unobserved fixed components.
- By the Frisch-Waugh-Lovell Theorem, this is equivalent to:
  - Regress  $Y$ , time FE,  $X$ , and  $D$  on full individual dummies and get residuals.
  - Regress  $Y$  residuals on  $X$ , time FE, and  $D$  residuals.
- Alternatively, use **First Differences** (FD):

$$Y_{it} - Y_{it-1} = \lambda_t - \lambda_{t-1} + (X'_{it} - X'_{it-1})\beta + \rho(D_{it} - D_{it-1}) + \varepsilon_{it} - \varepsilon_{it-1}$$

#### Key facts:

- Both FE and FD are unbiased estimators of the treatment effect.
- When  $T = 2$ , they yield the same result.
- With  $T > 2$ , efficiency differs depending on serial correlation in  $\varepsilon_{it}$ .
- with FD you reduce sample size by 50% in FE you lose dfs (1 df for each individual, N dfs lost! use `xtreg ... fe` to adjust standard errors for lost degrees of freedom.)

#### 1.4.3. Incidental Parameters Problem

- When  $N \rightarrow \infty$  and  $T$  is fixed, it is impossible to consistently estimate the  $N$  individual fixed effects as the number of parameters grows with  $N$ . Truly, the problem is solved just by large  $T$ .
- Incidental problem because it is emerging from thousands of parameters that are not what you want to identify but are necessary to identify
- When differencing from individual means, we avoid estimating fixed effects directly, thus solving the problem.

#### 1.4.4. FE and Measurement Error

- Fixed effects estimates are highly sensitive, more than OLS, to measurement error, which can lead to strong attenuation bias (maybe try IV).
- Proof:** Say that  $x = x^* + \nu$  and assume that:

$$\text{cov}(x_{it}^*, \varepsilon_{it}) = 0, \quad \text{cov}(\nu_{it}, \varepsilon_{it}) = 0$$

the first assumption rules out endogeneity between the true regressor (hence also the proxy as the proxy is the true plus white noise) and the error, the second assumption is that measurement error is uncorrelated over time (there is no autocorrelation in the measurement error).

#### OLS attenuation bias:<sup>1</sup>

<sup>1</sup>Proof: attenuation bias for OLS. Suppose the true model is:

$$y = \beta x^* + \varepsilon$$

but the true regressor  $x^*$  is not observed. Instead, we observe:

$$x = x^* + u$$

We estimate the regression:

$$y = \gamma x + v$$

and want to know whether  $\gamma = \beta$ . OLS gives:

$$\gamma = \frac{\text{Cov}(y, x)}{\text{Var}(x)}$$

You don't have the true model  $y = \beta x^* + \varepsilon$ , fill in  $x = x^* + u$ :

$$\gamma = \frac{\text{Cov}(\beta x^* + \varepsilon, x^* + u)}{\text{Var}(x^* + u)}$$

Compute the numerator:

$$\begin{aligned} \text{Cov}(\beta x^* + \varepsilon, x^* + u) &= \beta \text{Cov}(x^*, x^*) + \beta \text{Cov}(x^*, u) + \text{Cov}(\varepsilon, x^*) + \text{Cov}(\varepsilon, u) \\ &= \beta \text{Var}(x^*) \quad (\text{since all cross-covariances are zero}) \end{aligned}$$

$$\hat{\beta}_{OLS} = \beta \cdot \frac{\sigma_{x^*}^2}{\sigma_{x^*}^2 + \sigma_u^2} = \beta \left( 1 - \frac{\sigma_\nu^2}{\sigma_{x^*}^2 + \sigma_\nu^2} \right)$$

#### • Fixed Effects attenuation bias (with serial correlation):

$$\hat{\beta}_{FE} = \beta \left( 1 - \frac{\sigma_\nu^2}{(1 - \rho_x)\sigma_{x^*}^2 + \sigma_\nu^2} \right), \quad \rho_x = \frac{\text{cov}(x_{it}^*, x_{it-1}^*)}{\text{var}(x_{it}^*)}$$

- $\rho_x$  is the autocorrelation of the true regressor  $x_{it}^*$  across time for the same individual. When  $\rho_x = 0$ , FE and OLS yield the same bias. When  $\rho_x \approx 1$ , the difference  $x_{it}^* - x_{it-1}^*$  is small and the bias from measurement error increases. Why? As you can see  $\rho_x$  shrinks the true variability in  $\sigma_{x^*}^2$ , hence shrinks the signal to noise ratio (there is less signal, that is, information, around! you need more observations!).

## 2. Difference-in-Differences (DID)

### 2.1. Basic Setup

- DID is a special case with **two time periods** ( $t = 0$  pre-treatment,  $t = 1$  post-treatment) and two groups ( $g = 0$  control,  $g = 1$  treatment). We have a **within group variation** in treatment status (in the treated group).
- Does not require panel data, repeated cross-sections are sufficient.

### 2.2. Assumptions

Strong Assumption (which implies treatment effect is homogeneous and constant over time):

$$\mathbb{E}(Y_{0igt} | g, t) = \gamma_g + \lambda_t, \quad \mathbb{E}(Y_{1igt} | g, t) = \mathbb{E}(Y_{0igt} | g, t) + \rho$$

where:

- $\gamma_g$ : **group fixed effect** (the group version of  $\alpha_i$ , we will remove the time varying unobservables proper of the group!). Captures time-invariant differences between groups (e.g., region, demographics). Visually: the initial vertical gap between the lines.
- $\lambda_t$ : **time fixed effect**. Captures all time-varying shocks that affect all groups equally. Visually: a shift or slope change that applies equally to all groups over time.

### 2.2.1. Unconfoundedness and Parallel Trends

- Such a strong assumption **implies unconfoundedness, look at the two expected values above**, the treated outcome is just as the control outcome plus the treatment effect! That is, **In absence of treatment, outcomes would evolve similarly across groups, ie the pre-post change in the treatment group would have followed the same trend as the pre-post change in the control group if the treatment group had not been treated**:

$$\mathbb{E}[Y_{0it} - Y_{0it-1} | D_{it} = 1] = \mathbb{E}[Y_{0it} - Y_{0it-1} | D_{it} = 0]$$

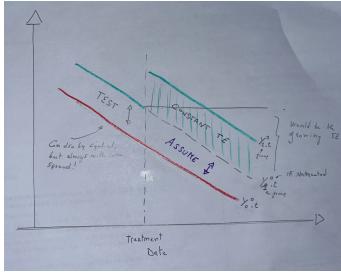
Note: the assumption above about parallel trends refers to both the pre-treatment trend (the famous pre-trend test) and the post-treatment trend (look at the  $t$  subscript). However, only the pre-trend can be tested empirically.

Compute the denominator:

$$\text{Var}(x^* + u) = \text{Var}(x^*) + \text{Var}(u)$$

Final expression:

$$\gamma = \beta \cdot \frac{\text{Var}(x^*)}{\text{Var}(x^*) + \text{Var}(u)} = \beta \cdot \frac{\sigma_{x^*}^2}{\sigma_{x^*}^2 + \sigma_u^2}$$



130 • Common violations of pre-trend:

- 131 – Only one group has a trend change due to external  
132 factor (eg another treatment policy). Violations can come  
133 from either group treatment or control.  
134 – Group composition changes over time. Violations can  
135 come from either group treatment or control.  
136 – **Ashenfelter dip:** pre-treatment drop for treated group.

137 • Robustness Checks:

- Visually: Graph trends before treatment for the two groups, check they are parallel (we need more than two periods). Note: If covariates affect the outcome trends, it may be better to condition on them (e.g., plot residualized trends after regressing out covariates).

$$\mathbb{E}[Y_{0it} - Y_{0it-1} | D_{it} = 1, X_{it}] = \mathbb{E}[Y_{0it} - Y_{0it-1} | D_{it} = 0, X_{it}]$$

- 138 – In the regression: Control for group-specific trends, re-  
139 sults should not change. This is because they can have dif-  
140 ferent pretrends, eg T is growing faster than C, and we need  
141 assume it would have continued to grow faster. more infra  
142 – Discussion: Add a discussion to Justify parallel trends as-  
143 sumption. Justify the absence of any other simultaneous  
144 treatment

145 **Remark 1:** Remark: The SUTVA is implicit in the unconfounded-  
146 ness assumption: a violation of the SUTVA would lead to an increase  
147 in the outcome of the control group after treatment (rhs of the equa-  
148 tion, see figure). In the plot, you have that with treatment not only  
149 the outcome of T shifts but also the outcome of control. Plus, to be  
150 precise, we would need to change our potential outcome framework!

151 **Remark 2** Under the parallel trends assumption, we are essen-  
152 tially saying that we want the treatment and control groups to be  
153 as similar as possible (they have to follow the same evolution). The  
154 goal is to ensure that assignment to treatment is not driven by in-  
155 trinsic characteristics, but rather that treatment was exogenously or  
156 randomly assigned (ideally due to some factor that cannot be ma-  
157 nipulated, such as random border assignment in Africa). However,  
158 the more similar the two groups are, the higher the risk of spillovers.  
159 Trade-off: The potential violation of SUTVA is the cost of ensuring  
160 that the treatment and control groups are highly comparable. There-  
161 fore, you must provide a convincing argument for why SUTVA holds  
162 in your setting.

164 **2.2.2. Regression Specification (Fully Saturated)**

165 Now we will use a fully saturated regression. Silly remark: in the re-  
166 gression above you had aggregated the constant in the group dummy.

- Model:

$$Y_{igt} = c_0 + c_1 D_g + c_2 D_t + c_3 D_g D_t + \varepsilon_{igt}$$

- We can estimate the treatment effect by taking the difference between the two differences in means:

$$\hat{\rho} = (\bar{Y}_{t=1,g=1} - \bar{Y}_{t=0,g=1}) - (\bar{Y}_{t=1,g=0} - \bar{Y}_{t=0,g=0}) \\ = (\text{Tr Post} - \text{Pre}) - (\text{Cr Post} - \text{Pre})$$

- Math:

$$\begin{aligned} & \mathbb{E}(Y_{igt} | g = 1, t = 1) - \mathbb{E}(Y_{igt} | g = 1, t = 0) \\ & - [\mathbb{E}(Y_{igt} | g = 0, t = 1) - \mathbb{E}(Y_{igt} | g = 0, t = 0)] \\ & = [(c_0 + c_1 + c_2 + c_3) - (c_0 + c_1)] - [(c_0 + c_2) - c_0] \\ & = (c_2 + c_3) - c_2 = c_3 \end{aligned}$$

- in the table below the horizontal rows remove the constant and the group effect (again, differencing removes unobserved components).
- the vertical line removes the time trend

	t=1	t=0	DIF
g=1	$c_0 + c_1 + c_2 + c_3$	$c_0 + c_1$	$c_2 + c_3$
g=0	$c_0 + c_2$	$c_0$	$c_2$
DIF	$c_1 + c_3$	$c_1$	$c_3$

Treatment effect:  $c_3$

171 **2.2.3. Regression specification**

The regression specification is:

$$Y_{igt} = \gamma_g + \lambda_t + \rho D_{gt} + \varepsilon_{igt}$$

173  $D_{gt} = 1$  only for  $g = 1$  and  $t = 1$ .

174 **2.2.4. Extensions**

- If you have **Repeated Cross-section:** you can pool the data and use (trivial)

$$Y_{igt} = c_0 + c_1 D_g + c_2 D_t + c_3 D_g D_t + \varepsilon_{igt}$$

- If you have **panel data.** Use FE, see below the regression at the individual level (with individual level partialling out)

$$Y_{it} = \alpha_i + \lambda_t + X'_{it}\beta + \rho D_{it} + \varepsilon_{it}$$

- you could also use the FD model

175 **3. Leads and Lags of Treatment (now we have multiple t)**

- Estimate dynamic effects:

$$Y_{igt} = \alpha_g + \lambda_t + \sum_{\tau=0}^m \rho_{-\tau} D_{g,t-\tau} + \sum_{\tau=1}^q \rho_{+\tau} D_{g,t+\tau} + X'_{igt}\beta + \varepsilon_{igt}$$

177 note  $\tau$  makes you go back in time!

- Purpose:

- Detect pre-trends (check if  $\rho_{+\tau} = 0$  before treatment)
- Observe treatment effect over time.

181 **3.1. Controlling for Covariates**

General Principle: We include covariates to avoid bias only when they affect the untreated potential outcome  $Y^0$ . VERY TRIVIAL:

- **Time-invariant or independent of treatment:** say 1) race is time-invariant, 2) educ varies but does not affect treatment

- If effect on outcome  $Y^0$  is constant: don't need to control (there is the individual level fe)
- If effect on outcome  $Y^0$  varies over time: include interaction with  $t$  (i.e.,  $t \times x$ ).

- **Time-varying and correlated with treatment:** say race changes and it changes with treatment

- If effect on the outcome  $Y^0$  is constant: include  $x_t$  in levels.
- If time-varying effect on the outcome  $Y^0$ : include  $t \times x_t$  interaction.

### 3.2. Group-Specific Trends

Suppose the treatment group was already growing faster than the control group before treatment. Then, any post-treatment difference in outcomes could be due to this ongoing trend, not the treatment itself. By adding group-specific trends to the regression, we are asking whether *there is still a treatment effect after accounting for the fact that the treatment group was already growing faster?*. Just make the plot; the gap is shrinking/expanding!<sup>2</sup>. Hence:

- Allow differential trends across groups (linear/polynomial):

$$Y_{igt} = \gamma_{0g} + \gamma_{1g}t + \lambda_t + X'_{igt}\beta + \rho D_{igt} + \varepsilon_{igt}$$

- Problem: we cannot control for non-linear trends differential across groups easily. Use triple differences!

### 3.3. Triple Differences

#### 3.3.1. Idea

- Motivation:** Difference-in-Differences (DID) may fail if group-specific trends are nonlinear or not parallel. Triple Differences (DDD) adds a third dimension to difference out these biases.
- Setup and Example:** Medicaid introduced in some states ( $g$ ), and only for families with young children ( $a$ ).

- **Setup:** We observe:

- \* Two groups ( $g$ ): a *treated region* and an *untreated region*.
- \* Two time periods ( $t$ ): before and after a policy intervention.
- \* Two subgroups ( $a$ ):
  - Treated subgroup:** families with children.
  - Untreated subgroup:** families without children.

- **Why Triple Differences (DDD) helps:**

- \* You would have compared families with child in the treated state with families without child in the treatment state.
- \* what if these two types of families intrinsically had different trends? (families without child are younger, and earning of young people are declining relatively more than earnings of reach people?)
- \* To correct for this differential trend, you partial out from those two groups the equivalent groups in another state.

#### 3.3.2. Regression Specification:

$$Y_{iagt} = \gamma_{gt} + \lambda_{at} + \theta_{ag} + \rho D_{agt} + \varepsilon_{iagt}$$

where:

- $D_{agt}$  is a dummy for treated group (state  $g$ )  $\times$  treated subgroup  $a$  (families with child)  $\times$  post period.
- $\gamma_{gt}$ : group-time fixed effects (controls for group-specific trends).
- $\lambda_{at}$ : subgroup-time fixed effects.
- $\theta_{ag}$ : group-subgroup fixed effects (e.g., in the treated state, families without children may systematically earn less). being in that subgroup in that state creates an advantage per se, baseline.

#### Alternative Saturated Specification

$$Y_{iagt} = c_0 + c_1 D_g + c_2 D_t + c_3 D_g D_t + c_4 D_a + c_5 D_t D_a + c_6 D_g D_a + c_7 D_g D_t D_a + \varepsilon_{iagt}$$

- $c_7$  identifies the DDD estimate

<sup>2</sup>Note that this formula no longer works as it is based on the parallel trend assumption

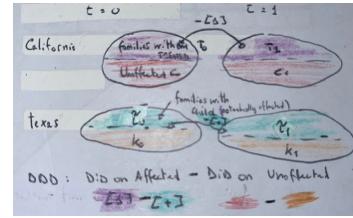
$\hat{\rho} = (\bar{Y}_{t=1,g=1} - \bar{Y}_{t=0,g=1}) - (\bar{Y}_{t=1,g=0} - \bar{Y}_{t=0,g=0}) = (\text{Treated Post - Pre}) - (\text{Control Post - Pre})$

### 3.3.3. Non-Parametric Specification

#### Color Legend:

- Treated region, families with children (Group A)
- Treated region, families without children (Group B)
- Untreated region, families with children (Group A)
- Untreated region, families without children (Group B)

$$\mathbb{E}(Y_{igt} | g=1, t=1, a=1) - \mathbb{E}(Y_{igt} | g=1, t=0, a=1) - \mathbb{E}(Y_{igt} | g=0, t=1, a=1) + \mathbb{E}(Y_{igt} | g=0, t=0, a=1) \\ - [\mathbb{E}(Y_{igt} | g=1, t=1, a=0) - \mathbb{E}(Y_{igt} | g=1, t=0, a=0) - \mathbb{E}(Y_{igt} | g=0, t=1, a=0) + \mathbb{E}(Y_{igt} | g=0, t=0, a=0)]$$



- you can control both for the differential trends between states (obvious baseline) and the differential trends between groups.

### 3.4. Standard Errors

- Problem:** DID estimates can severely underestimate standard errors due to serial correlation.

#### Background:

- Downwards bias arising from time series structure in the data (serial correlation).
- Persistency reduces the amount of information contained in the data (low signal to noise ratio). Also, I guess, lower effective sample size.

#### Solutions and Alternatives:

- Collapse data into pre- and post-treatment periods:** Aggregating eliminates serial correlation within units over time, making standard errors more reliable.
- Specify the correlation structure directly:** Use autoregressive models (e.g., AR(1)) for the error term.
- Block bootstrap:** Resample at the group level (e.g., states or clusters). Valid when the number of blocks is large.
- Clustered Standard Errors:** Allows for arbitrary serial correlation within each group (e.g., state, firm). Consistent only when the number of clusters is large.

### 3.5. DID with Treatment at Different Times (Staggered Design)

#### 3.5.1. Framing the Issue

##### • Canonical Model (TWFE):

$$Y_{igt} = \gamma_g + \lambda_t + \rho D_{gt} + \varepsilon_{igt}$$

- $\gamma_g$ : group fixed effects
- $\lambda_t$ : time fixed effects
- $D_{gt}$ : treatment indicator for group  $g$  at time  $t$
- $\rho$ : estimated average treatment effect

##### • Interpretation of $\rho$ in staggered adoption:

$$\rho_{fe} = \mathbb{E} \left( \sum_{(g,t): D_{gt}=1} W_{g,t} \Delta_{g,t} \right)$$

where:

- $\Delta_{g,t}$ : treatment effect in group  $g$  at time  $t$
- $W_{g,t}$ : weight assigned to cell  $(g, t)$ , with  $\sum W_{g,t} = 1$

- Problem:** some  $W_{g,t}$  may be negative.

- 283 - Even if all  $\Delta_{g,t} > 0$ , negative weights can make  $\rho_{fe} < 0$ .  
 284 - Arises due to comparisons across differentially treated  
 285 groups in TWFE.  
 \* note treatment timing may be endogenous: groups  
 with higher treatment benefit are treated before

- 286 • **Weight decomposition (de Chaisemartin):** We will build  
 287 the ATE as a weighted average of the two by two comparisons  
 288 in treatment effects. The weights in such weighted average are  
 289 determined by the leverage of the treatment dummy on the re-  
 290 gression:

$$D_{gt} = \alpha + \gamma_g + \lambda_t + \varepsilon_{gt} \Rightarrow \varepsilon_{gt} = D_{gt} - \alpha - \gamma_g - \lambda_t = D_{gt} - \underbrace{D_{.,t}}_{D^{pred}}$$

We need the above because the weights are driven by the errors  $\varepsilon_{gt}$ ! These residuals measure how much the treatment dummy  $D_{gt}$  deviates from what would be expected based on group and time averages.

$$\rho_{FE} = \mathbb{E} \left( \sum_{(g,t): D_{gt}=1} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t} \right), \quad \text{where } w_{g,t} = \frac{\varepsilon_{gt}}{\sum_{(g,t): D_{gt}=1} \frac{N_{g,t}}{N_1} \varepsilon_{gt}}$$

- 291 – so if  $\varepsilon_{gt}$  is negative the weight is negative!  
 292 – Later periods of early adopters are very exposed to nega-  
 293 tive weights because: 1) the group is usually treated, as  
 294 it was treated from earlier stages, and 2) in later stages  
 295 many groups have been treated. Early adopters in late pe-  
 296 riods face high expected treatment! *My Idea:* The weight  
 297 is negative if  $D - D^{pred}$  is negative, so if the predicted level  
 298 of treatment is greater than the actual one that is, if we  
 299 are greatly expecting treatment (consider that the linear  
 300 regression does not know that there is a 01 bound, so it is  
 301 possibly giving to an early adopter a predicted value  $> 1$ ).  
 302 – With homogeneous treatment effects, you still have nega-  
 303 tive weights, but since the weights sum to one, negative  
 304 weights are compensated by positive ones, and both pos-  
 305 itive and negative weights are multiplying the same TE,  
 306 so they cancel out. The issue is that if those groups with  
 307 negative weights have higher treatment effects, the ATE is  
 308 negative!

• **Illustrative example:**

- 309 – Group 1: treated only in period 3  
 310 – Group 2: treated in periods 2 and 3  
 311 – Same group sizes  
 – **Step-by-step computation of residuals**  $\varepsilon_{gt} = D_{gt} - D_{.,t} + D_{.,.}$

\* **Step 1: Treatment status  $D_{gt}$**

	$t = 1$	$t = 2$	$t = 3$
$g = 1$	0	0	1
$g = 2$	0	1	1

The dummies represent the probability of being treated

\* **Step 2: Row averages:**  $D_{1,.}$  represents the probability of being treated of group 1 across all  $t$ .

$$D_{1,.} = \frac{0+0+1}{3} = \frac{1}{3}, \quad D_{2,.} = \frac{0+1+1}{3} = \frac{2}{3}$$

\* **Step 3: Column averages:**  $D_{.,1}$  represents the probability of being treated in time  $t$

$$D_{.,1} = \frac{0+0}{2} = 0, \quad D_{.,2} = \frac{0+1}{2} = \frac{1}{2}, \quad D_{.,3} = \frac{1+1}{2} = 1$$

\* **Step 4: Grand average**

$$D_{.,.} = \frac{0+0+1+0+1+1}{6} = \frac{3}{6} = \frac{1}{2}$$

\* **Step 5: Compute residuals**

$$\begin{aligned} \varepsilon_{1,3} &= D_{1,3} - D_{1,.} - D_{.,3} + D_{.,.} = 1 - \frac{1}{3} - 1 + \frac{1}{2} = \frac{1}{6} \\ \varepsilon_{2,2} &= D_{2,2} - D_{2,.} - D_{.,2} + D_{.,.} = 1 - \frac{2}{3} - \frac{1}{2} + \frac{1}{2} = \frac{1}{3} \\ \varepsilon_{2,3} &= D_{2,3} - D_{2,.} - D_{.,3} + D_{.,.} = 1 - \frac{2}{3} - 1 + \frac{1}{2} = -\frac{1}{6} \end{aligned}$$

–  $\varepsilon_{1,3} = 1/6, \varepsilon_{2,1} = 1/3, \varepsilon_{2,3} = -1/6$

– Then:

$$\rho_{fe} = \frac{1}{2} \mathbb{E}[\Delta_{1,3}] + \mathbb{E}[\Delta_{2,2}] - \frac{1}{2} \mathbb{E}[\Delta_{2,3}]$$

– Suppose:

$$\mathbb{E}[\Delta_{1,3}] = \mathbb{E}[\Delta_{2,2}] = 1, \quad \mathbb{E}[\Delta_{2,3}] = 4 \Rightarrow \rho_{fe} = \frac{1}{2}(1) + 1 - \frac{1}{2}(4) = -0.5$$

– Despite all ATEs being positive, the estimated  $\rho_{fe}$  is negative.

This is connected to the cross-comparison interpretation of the problem (beacon): The issue arises when later-treated groups are compared to earlier-treated groups, and the earlier-treated group is incorrectly used as a control, even though it's already treated. So if you have 1) a group that has been treated for long, this will be a poor control for the other recently treated compared units and will create a negative ATE, 2) if you have many treated groups at  $t$ , then it is more likely to get a poor comparison.

the fun part is that these poor counterfactuals are likely exactly those with high ate. they are terrible counterfactuals and they are drawing down the ate!

### 3.6. Solution

- how to notice there is a problem? compute standard deviation of the ATEs across the treated (g,t) cell (where negative weights may come in) and compare it with the absolute value of the DID estimate (e.g., using not yet treated as control!).
- Solution:** Estimate treatment effects by comparing treated units only to *not-yet-treated* units. Similar to matching vs OLS: drop poor controls!

## 4. DID Advanced Topics

### 4.1. Lagged Dependent Variable

- Sometimes we **must include the lagged outcome**  $Y_{t-1}$  as a covariate (control for it).
- Why include the lagged outcome ( $Y_{t-1}$ )? Ashenfelter's Dip and Pre-trends Bias**

If treatment  $D_t$  depends on the past outcome  $Y_{t-1}$ , and  $Y_{t-1}$  is correlated with the current outcome  $Y_t$ , then  $D_t$  is endogenous: it is correlated with  $Y_t$  through  $Y_{t-1}$ . This leads to omitted variable bias if  $Y_{t-1}$  is not controlled for. Note: this is, first and foremost, a violation of parallel trends. Also, recognize that this is selection bias: selection bias is what causes the violation of parallel trends. This is self-selection because the treated group shares something in common that makes them more likely to enter treatment. Most likely, they were chosen to be treated **because** they experienced a drop in income. There is no internal validity: pre-trends are violated, and the treatment group is not comparable to the control group.

\*Be careful not to misunderstand: you may think "had the control group not received treatment, their  $Y_0$  would have been the same."

357      *No: the control group experienced a drop in income, and that's why they were selected. This clearly creates positive bias!*  
 358  
 359  
 360      *Graph intuition:* Initially, treated and control groups have parallel trends. A dip in  $Y_{t-1}$  for treated units breaks this trend. You want to control for  $Y_{t-1}$  that is, absorb the dip. Including  $Y_{t-1}$  restores the trend. Imagine the "smoothing" effect a time dummy would have: a time dummy is common to all individuals, but  $Y_{t-1}$  provides an individual-specific intensity.  
 361  
 362  
 363  
 364  
 365      •  $Y_{t-1}$  is clearly not eliminated by individual fixed effects.  
 366      • In OLS, including  $Y_{t-1}$  is straightforward, but in Fixed Effects (FE) models, this becomes more complex (see subsection).

#### 368      4.1.1. OLS with Lagged Dependent Variable

- One possibility is to **assume unconfoundedness** given lagged outcomes even without controlling for fixed effects (basically say you don't need FE):

$$\mathbb{E}[Y_{0it} \mid Y_{it-k}, X_{it}, D_{it}] = \mathbb{E}[Y_{0it} \mid Y_{it-k}, X_{it}]$$

- Estimate via OLS or Random Effects (RE)

$$Y_{it} = \alpha + \theta Y_{it-k} + \lambda_t + \rho D_{it} + X'_{it} \beta + \varepsilon_{it}$$

369      no nickel bias (more infra) in OLS adn RE (no differencing  
 370      trasformation)

#### 371      4.2. Lagged Dependent Variable with Fixed Effects

- When we assume unconfoundedness given lagged outcomes and fixed effects:

$$\mathbb{E}[Y_{0it} \mid Y_{it-k}, \alpha_i, X_{it}, D_{it}] = \mathbb{E}[Y_{0it} \mid Y_{it-k}, \alpha_i, X_{it}]$$

372      the treatment status is independent of the potential outcome.

- Estimation equation:

$$Y_{it} = \alpha_i + \theta Y_{it-k} + \lambda_t + \rho D_{it} + X'_{it} \beta + \varepsilon_{it}$$

- Problem: **Nickell bias** (see below).

- Consider first-differencing (fixed effects disappear):

$$\Delta Y_{it} = \theta \Delta Y_{it-1} + \Delta \lambda_t + \rho \Delta D_{it} + \Delta X'_{it} \beta + \Delta \varepsilon_{it}$$

- Problem!  $\Delta \varepsilon_{it} = \varepsilon_{it} - \varepsilon_{it-1}$  is correlated with  $\Delta Y_{it-1}$ .
- $\Delta Y_{it-1}$  contains  $\varepsilon_{it-1}$ , which is in  $\Delta \varepsilon_{it}$ .
- $\Rightarrow$  correlation between regressor and error term  $\Rightarrow$  bias.
- Nickell (1981) proved that OLS is inconsistent in short panels.
- Solution:

- \* If  $\varepsilon_{it}$  is serially correlated ( $\text{corr}(\varepsilon_{it}, \varepsilon_{it-k}) \neq 0$ ). Then there might not be any solution.
- \* If the correlation does not hold up to period k: use  $Y_{it-2-k}$  as IV for  $Y_{it-1}$  (or First-differences:  $Y_{it-2} - Y_{it-3}, X_{it-2} - X_{it-3}$ , etc. valid instrument for an AR(1) process). They become uncorrelated (trivial). if there is maximum AR(k) correlation the instrument is valid, but we **will also need to check the first stage** (not trivial).
- \* Empirically we can bound the effect of interest:

- If we include fixed effects and no lagged outcome, we expect upwards biased estimates (we have a dip and then up: Dip in  $Y_{it-1}$  leads to mechanical rebound; FE attributes this to treatment  $\rightarrow$  overestimates effect).
- If we do not include fixed effects but we include lagged outcomes, we expect downwards biased estimates (boh)

#### 5. Staggered DID Rho interpretation

398      Weighted average of all possible two-by-two DID estimators. Only  
 399      recovers the ATE when treatment effects are homogeneous  
 400

#### 401      6. DID in levels vs DID in logs

- **Key assumption: common trends**, which is *not* invariant to monotonic transformations of the outcome (parallel trend in levels are not parallel trends in logs). Thsi is obvious bc in levels you have the difference between the outcomes t obe fixed, in logs you have the ratio to be fixed
- **Treatment Effects found using DID can differ greatly depending on whether levels or logs are used** as the outcome variable.