

Jorge Eduardo Pérez Pérez

# Advances in Differences in Differences and Bartik instruments: Class Notes

CEER, Banco de la República de Colombia - Cartagena

October 9, 2023



## Chapter 2

# Staggered adoption, heterogeneity, and issues with the TWFE specification

In this chapter we introduce the staggered adoption setup and discuss issues with the TWFE specification for estimation of treatment effects under staggered adoption. We follow the setup in Roth et al. (2023).

### 2.1 Staggered adoption setup

The staggered adoption setup is motivated by units selectively being treated over time. For example, we can think of units as cities and treatment being the entry of a ride-sharing platform such as Uber to those cities. Uber does not enter all cities at the same time, and once Uber enters a city it does not leave (or rarely does).

Now we assume there are  $T$  periods indexed by  $t = 1, \dots, T$ . The variable  $D_{i,t}$  denotes treatment status for unit  $i$  at time  $t$ . For staggered adoption, we assume the following:

**Assumption 2.1** (*Staggered adoption*)

1. *Treatment is binary:*  $D_{i,t} \in \{0, 1\}$ .
2. *All units begin untreated:*  $D_{i,1} = 0$  for all  $i$ .
3. *Treatment is absorbing:*  $D_{i,t'} \geq D_{i,t}$  for all  $i$  and  $t' \geq t$ .

The absorbing treatment assumption is convenient because it lets us group the treated units by treatment cohorts. We define the treatment cohort  $G_i$  as  $G_i = \min \{t : D_{i,t} = 1\}$ , the first period when an unit receives treatment. For control units (that never receive treatment),  $G_i = \infty$ .

With multiple time periods, we also need to extend the potential outcomes notation to allow for treatment histories. A treatment history is a vector of length  $T$  of zeros and ones. Moreover, in the staggered adoption setting, all the treatment history vectors will be of the form  $(\mathbf{0}_s, \mathbf{1}_{T-s})$ .

The potential outcome for unit  $i$  if they were treated for the first time at time  $g$  (that is, if their cohort  $G_i = g$ ) is  $Y_{i,t}(\mathbf{0}_{g-1}, \mathbf{1}_{T-g+})$ . If it was never treated, its potential outcome is  $Y_{i,t}(\mathbf{0}_T)$ . We can simplify these in the staggered adoption setting as  $Y_{i,t}(g) = Y_{i,t}(\mathbf{0}_{g-1}, \mathbf{1}_{T-g+})$  and  $Y_{i,t}(\infty) = Y_{i,t}(\mathbf{0}_T)$ .

We can also extend our notation for treatment effects. The **individual treatment effect** for unit  $i$  if they belong to cohort  $g$  is

$$\tau_{i,t}(g) = Y_{i,t}(g) - Y_{i,t}(\infty). \quad (2.1)$$

Notice subtle differences with the treatment effect in the simple case. This is defined in terms of potential outcomes.

With the individual treatment effects as building blocks, you can build other quantities of interest. We will focus on the **average treatment effect on the treated for cohort  $g$  at time  $t$**  ( $ATT_{g,t}$ ) is the average of these individual treatment effects for a particular treatment cohort:

$$\tau_{g,t} = \mathbb{E}[\tau_{i,t}(g)|G_i = g]. \quad (2.2)$$

You could also define an **average treatment effect on the treated at time t**:

$$\tau_t = \mathbb{E}[\mathbb{E}[\tau_{i,t}(g)|G_i = g]].$$

To identify  $ATT_{g,t}$ , we need to generalize the parallel trends assumption. In the basic case, we assumed that in absence of treatment, the outcomes of the treated and control groups would evolve in parallel. A natural extension is to require that the untreated potential outcomes of all cohorts evolve in parallel:

**Assumption 2.2** (*Strong unconditional parallel trends*) For all  $t \neq t'$  and  $g \neq g'$ :

$$\mathbb{E}[Y_{i,t}(\infty) - Y_{i,t'}(\infty)|G_i = g] = \mathbb{E}[Y_{i,t}(\infty) - Y_{i,t'}(\infty)|G_i = g']$$

We also extend the no anticipation assumption as

**Assumption 2.3** (*Staggered no anticipation*)

$$Y_{i,t}(g) = Y_{i,t}(\infty) \text{ for all } i \text{ and } t < g$$

## 2.2 Two-way fixed effects estimation

To introduce two-way fixed effects estimation, assume that treatment effects are constant across time and across units:

$$\tau_{g,t} = \tau \text{ for all } t \geq g$$

In this scenario, a natural extension of (??) is a linear model with unit and time effects:

$$Y_{i,t} = \alpha_i + \theta_t + D_{i,t}\beta + \varepsilon_{i,t} \quad (2.3)$$

Under constant treatment effects,  $\hat{\beta}$  from (2.3) is a consistent estimator for  $\tau$ , and it can be estimated through fixed-effects estimation of (2.3) (For review of fixed effects estimation, see ? ).

## 2.3 Treatment effect heterogeneity and TWFE weights

If we are willing to assume that treatment effects are constant, there is nothing wrong with TWFE estimation. There may be settings where this is a plausible assumption: for example, all units are treated at the same time and economic theory suggests that there is not treatment heterogeneity. Or, units are treated at different times but treatments have a one-off, homogeneous impact.

However, in most applications, we may expect **heterogeneity in treatment effects**. Economic theory will often imply heterogeneous treatment effects across units. For example, the elasticity of labor supply will be different across individuals with different outside options. Even in settings where the per-period treatment effect is homogeneous, we may expect differences at time  $t$  across units that belong to different cohorts. For example, the effects of Uber entry in a city on traffic congestion may be different after two years vs. after one year, as the size of Uber's fleet grows.

We would expect that the estimate from (2.3) corresponds to a weighted average of  $ATT_{g,t}$  with some reasonable weights, i.e.,  $\hat{\beta} = \sum_{t,g} \omega_{t,g} \tau_{t,g}$  with  $\omega_{t,g} > 0$ . However, this turns out not to be the case. To show some mathematical intuition of this, note that by the FWL theorem, the OLS estimate of  $\beta$  from

(2.3) equals the coefficient of a regression of  $Y_{it}$  on the residuals of a regression of  $D_{it}$  on unit and time effects. That is, run the regression  $D_{it} = \tilde{\alpha}_t + \tilde{\delta}_i + u_{it}$  and obtain the residuals  $D_{it} - \hat{D}_{it}$ . Then the OLS estimate of  $\beta$  from (2.3) equals:

$$\hat{\beta} = \frac{\text{Cov}(Y_{it}, \hat{D}_{it} - D_{it})}{\text{Var}(D_{it} - \hat{D}_{it})} = \frac{\sum_{i,t} Y_{it} (D_{it} - \hat{D}_{it})}{\sum_{i,t} (D_{it} - \hat{D}_{it})^2} \quad (2.4)$$

If we break down the numerator into observations where  $D_{it} = 1$  and  $D_{it} = 0$  we can write:

$$\hat{\beta} = \frac{\sum_{i,t, D_{it}=0} Y_{it} (-\hat{D}_{it}) + \sum_{i,t, D_{it}=1} Y_{it} (1 - \hat{D}_{it})}{\sum_{i,t} (D_{it} - \hat{D}_{it})^2}$$

Moreover, when  $D_{it} = 1$ ,  $\tau_{i,t}(g) = Y_{i,t}(g) - Y_{i,t}(\infty) = Y_{i,t} - Y_{i,t}(\infty)$ . Replacing this value of  $Y_{i,t}$  in the numerator for  $D_{it} = 1$ :

$$\hat{\beta} = \frac{\sum_{i,t, D_{it}=0} Y_{i,t} (-\hat{D}_{it}) + \sum_{i,t, D_{it}=1} (Y_{i,t}(\infty) + \tau_{i,t}(g))(1 - \hat{D}_{it})}{\sum_{i,t} (D_{it} - \hat{D}_{it})^2} \quad (2.5)$$

Here, we can see that the OLS estimate of  $\beta$  equals a weighted average of treated and control observations. For treated observations, the weights are proportional to  $1 - \hat{D}_{it}$ . However, since  $\hat{D}_{it}$  is a prediction from a linear model, nothing guarantees that  $1 - \hat{D}_{it}$  will be positive! So some of the treatment effects  $\tau_{i,t}(g)$  will get negative weights when building  $\hat{\beta}$ , which is counterintuitive. Moreover, the weights in (2.5) need not be proportional to the sample sizes of each cohort  $g$ .

The negative weights will tend to arise for early-treated units, that is, units with low value of  $g$  late in the sample. The predicted value  $\hat{D}_{it}$  equals  $\bar{D}_i + \bar{D}_t - \bar{D}$ , where the bars denote sample means. Early treated units will have high values of  $\bar{D}_i$  (because they are treated in almost every  $t$ ) so  $\bar{D}_i \approx 1$ . If they are late in the sample then most units will have been treated, so  $\bar{D}_t \approx 1$ . Then  $\hat{D}_{it} \approx 2 - \bar{D}$ . The average  $\bar{D}$  will be less than one if there are non-treated units. So  $\hat{D}$  will be strictly higher than 1, and  $1 - \hat{D}_{it}$  will be negative in those cases.

## 2.4 Goodman-Bacon Decomposition

Goodman-Bacon (2021) proposes an intuitive decomposition of the TWFE estimator that illustrates why it may give counterintuitive weights to some units and how it involves “forbidden” comparisons that may lead the TWFE estimate to be a biased estimate of the average treatment effect.

Suppose there are three treatment cohorts: an “early-treated” cohorts with  $G_i = k$ , a “late-treated” treatment cohort with  $G_i = \ell > k$ , and an untreated cohort with  $G_i = \infty$ . We will also denote this cohort by  $U$  for untreated. We call  $PRE(k)$  as the time window before the  $k$  cohort is treated,  $MID(k, \ell)$  as the time window when the  $k$  cohort has been treated but the  $\ell$  cohort has not, and a  $POST(\ell)$  window when both cohorts have been treated. We denote by  $\bar{Y}_k^{PRE(k)}$  as the sample average of the outcome for units in the  $k$  cohort during the  $PRE(k)$  time window, and define other sample averages  $\bar{Y}_g^W, W \in \{PRE(k), MID(k, \ell), POST(\ell)\}, g \in \{k, \ell, \infty\}$  accordingly.

Let the fractions of units belonging to each cohort be  $n_k, n_\ell$ , and  $n_U$ , respectively. Assume that we observe a balanced panel.

[GB Fig 1. o Cunningham II 3 65]

With these three groups we can define four 2x2 difference in difference estimators:

[Cunningham II - 66-69]

A. Early treated vs. untreated:

$$\hat{\beta}_{kU}^{2 \times 2} = \left( \bar{Y}_k^{POST(k)} - \bar{Y}_k^{PRE(k)} \right) - \left( \bar{Y}_U^{POST(k)} - \bar{Y}_U^{PRE(k)} \right) \quad (2.6)$$

This comparison uses a fraction  $n_k + n_U$  of units, and since the panel is balanced, it also uses a fraction  $n_k + n_U$  of the  $NT$  observations in the panel.

B. Late treated vs. untreated:

$$\hat{\beta}_{\ell U}^{2 \times 2} = \left( \bar{Y}_\ell^{POST(\ell)} - \bar{Y}_\ell^{PRE(\ell)} \right) - \left( \bar{Y}_U^{POST(\ell)} - \bar{Y}_U^{PRE(\ell)} \right) \quad (2.7)$$

This comparison uses a fraction  $n_{\ell U} + n_U$  of units, and since the panel is balanced, it also uses a fraction  $n_{\ell U} + n_U$  of the  $NT$  observations in the panel.

C. Early treated vs. late treated:

$$\hat{\beta}_{k\ell}^{2 \times 2} = \left( \bar{Y}_k^{MID(k,\ell)} - \bar{Y}_k^{PRE(k)} \right) - \left( \bar{Y}_\ell^{MID(k,\ell)} - \bar{Y}_\ell^{PRE(\ell)} \right) \quad (2.8)$$

This comparison uses a fraction  $n_k + n_\ell$  of units. It does not use all the time periods, though: it only uses the periods in the  $PRE(\ell)$  window. Letting  $\bar{D}_k$  and  $D_\ell$  denote the fraction of periods in which each cohort is treated, this comparison uses a fraction  $(n_k + n_\ell)(1 - \bar{D}_\ell)$  of the  $NT$  observations.

D. Late treated vs. early treated:

$$\hat{\beta}_{\ell k}^{2 \times 2} = \left( \bar{Y}_\ell^{POST(\ell)} - \bar{Y}_\ell^{MID(k,\ell)} \right) - \left( \bar{Y}_k^{POST(\ell)} - \bar{Y}_k^{MID(k,\ell)} \right) \quad (2.9)$$

This comparison uses a fraction  $n_k + n_\ell$  of units. It does not use all the time periods, though: it only uses the periods in the  $POST(k)$  window. Letting  $\bar{D}_k$  and  $D_\ell$  denote the fraction of periods in which each cohort is treated, this comparison uses a fraction  $(n_k + n_\ell)(\bar{D}_k)$  of the  $NT$  observations.

Goodman Bacon shows that the TWFE estimate is a weighted average of these four difference-in-difference estimates. To understand the weights, recall that the coefficient estimate of linear regression on a binary variable and covariates (e.g. on a treatment indicator) will put more weight on covariate cells where there is most variation in the treatment. (MHE 3.3) Here, the weights will be proportional to the variance of treatment in each one of the comparisons, after adjusting for unit and time effects.

The variance of treatment for each one of the comparisons is:

$$\begin{aligned} \hat{V}_{jU}^D &= n_{jU}(1 - n_{jU})\bar{D}_j(1 - \bar{D}_j), j = k, \ell \\ \hat{V}_{k\ell}^D &= n_{k\ell}(1 - n_{k\ell})\frac{\bar{D}_k - \bar{D}_\ell}{1 - \bar{D}_\ell}\frac{1 - \bar{D}_k}{1 - \bar{D}_\ell} \\ \hat{V}_{\ell k}^D &= \frac{\bar{D}_\ell}{\bar{D}_k}\frac{\bar{D}_k - \bar{D}_\ell}{\bar{D}_k} \end{aligned}$$

where  $n_{ab} \equiv \frac{n_a}{n_a + n_b}$ .

With these variances, we can write the decomposition of the TWFE estimate:

$$\hat{\beta} = \sum_{k \neq U} s_{kU} \hat{\beta}_{kU}^{2 \times 2} + \sum_{k \neq U} \sum_{\ell > k} \left[ s_{k\ell} \hat{\beta}_{k\ell}^{2 \times 2} + s_{\ell k} \hat{\beta}_{\ell k}^{2 \times 2} \right] \quad (2.10)$$

with weights proportional to the variance of treatment and the sample size in each comparison:

$$\begin{aligned}
s_{kU} &= \frac{(n_k + n_U)^2 \hat{V}_{kU}^D}{\hat{V}^D}, \\
s_{k\ell} &= \frac{((n_k + n_\ell)(1 - \bar{D}_\ell))^2 \hat{V}_{k\ell}^D}{\hat{V}^D}, \\
s_{\ell k} &= \frac{((n_k + n_\ell)(\bar{D}_k))^2 \hat{V}_{\ell k}^D}{\hat{V}^D}.
\end{aligned}$$

The weights tend to be higher for comparisons where there’s more variance treatment, that is, groups where treatment occurs in the middle of the panel. This may be undesirable if we want to estimate ATT, because it will downweight some groups.

This decomposition tells us about what TWFE estimates, but it does not tell us whether it is unbiased for ATT or not. Let’s assume dynamic treatment effects as earlier and  $ATT_k(W)$  denote the average treatment effect in a treatment window for group  $k$ , e.g.  $ATT_k(W) = \frac{1}{T_W} \sum_{t \in W} \mathbb{E}[Y_{it}(k) - Y_{it}(0)]$ .

In this case each of the 2 by 2 estimates converges in probability to a different quantity:

$$\begin{aligned}
\beta_{kU}^{2 \times 2} &= ATT_k(Post) + \Delta Y_k^0(Post(k), Pre(k)) - \Delta Y_U^0(Post(k), Pre) \\
\beta_{k\ell}^{2 \times 2} &= ATT_k(MID) + \Delta Y_k^0(MID, Pre) - \Delta Y_\ell^0(MID, Pre) \\
\beta_{\ell k}^{2 \times 2} &= ATT_{\ell, Post(\ell)} + \Delta Y_{\ell\ell}^0(Post(\ell), MID) - \Delta Y_k^0(Post(\ell), MID) - (ATT_k(Post) - ATT_k(Mid))
\end{aligned}$$

The two first comparisons are “unproblematic”. Under parallel trends, these comparisons converge to ATTs for the given window. The third comparison, though, the later treated vs. early treated, is problematic under dynamic treatment effects. The difference-in-differences comparisons using the early treated units as controls is a “forbidden comparison”, because under dynamic treatment effects, the early treated may still have treatment effects happening when they are used as a control group for the later-treated groups. This is the same point risen by Borusyak and Jaravel (2018). In extreme settings, the contamination in this comparison may even flip the sign of the TWFE estimator!

Bacon’s decomposition has become a standard diagnostic tool to assess whether this contamination may be a potential issues. The decomposition consists of calculating the weights in (2.10) and see if there is a large weight on the late treated vs. early treated comparisons. Large weights on these comparisons are suggestive of potential issues in a scenario with dynamic treatment effects.

## 2.5 Callaway and Sant’Anna’s estimator for difference in differences with multiple treatment periods

Callaway and Sant’Anna (2021) take a different approach to estimating  $ATT_{g,t}$  and  $ATT$ . Since TWFE estimators may suffer from bias to estimate the ATT when there is staggered treatment adoption, and this bias may appear because TWFE makes forbidden comparisons, why not estimate all the  $ATT_{g,t}$ ’s separately and then aggregate them however we like? We can choose to use only the comparisons we like for aggregation. CS argue that by making this switch, we can focus on identification and easy to interpret estimates at the cost of harder implementation, instead of focusing on easy-to-implement but hard to interpret estimates such as those coming from TWFE.

CS’s assumptions are the same as those in the beginning of the section, although they do allow for weaker versions of the no anticipation and parallel trends assumptions. They also consider “conditional” versions of their assumptions instead of unconditional versions.

**Assumption 2.4** (*Limited anticipation*) There is a known  $\delta \geq 0$  such that

$$\mathbb{E}[Y_{i,t}(g)|X, G_i = g] = \mathbb{E}[Y_{i,t}(\infty)|X, G_i = g] \text{ for all } g \text{ such that } t < g - \delta$$

This allows for anticipation effects of the treatment up to  $\delta$  periods before the treatment. For parallel trends, we can choose either of two assumptions:

**Assumption 2.5** (Conditional parallel trends based on an untreated group) For all  $t$ , for each  $g$  such that  $t \geq g - \delta$

$$\mathbb{E}[Y_{it}(\infty) - Y_{i,t-1}(\infty)|X, G_i = g] = \mathbb{E}[Y_{i,t}(\infty) - Y_{i,t-1}(\infty)|X, G_i = \infty]$$

**Assumption 2.6** (Conditional parallel trends based on a not-yet-treated group) For all  $t$ , for each  $g$  such that  $t \geq g - \delta$

$$\mathbb{E}[Y_{it}(\infty) - Y_{i,t-1}(\infty)|X, G_i = g] = \mathbb{E}[Y_{i,t}(\infty) - Y_{i,t-1}(\infty)|X, G_i > g]$$

One additional assumption that the CS estimator needs is overlap. This is needed because it uses propensity score methods. Define the generalized propensity score as

$$p_{g,t}(X) = \Pr[G_i = g|X, \mathbf{1}(G_i = g) + C = 1] \quad (2.11)$$

where  $C$  is a dummy variable for  $G_i = \infty$ . This is the probability of being in treatment cohort  $g$  conditional on covariates and on being a member of group  $g$  or being in the control group. The overlap assumption is:

**Assumption 2.7** (Overlap) There exists  $\varepsilon > 0$  such that  $\Pr(G_i = g) > \varepsilon$  and  $p_{g,t}(X) < 1 - \varepsilon$

This means that for every covariate bin there are units in treatment cohort  $g$  and in the control cohort.

Before introducing CS's estimator we need to recall a few results from regression and propensity score matching (see Angrist and Pischke (2009)) to understand the inverse probability weighting estimator.

Think of a binary treatment  $D_i = 0, 1$  and potential outcomes  $Y_i(0), Y_i(1)$ . We want to estimate the ATT  $\mathbb{E}[Y_i(1) - Y_i(0)|D_i = 1]$ . The comparison between averages for treated and untreated does not identify ATT because of selection bias:

$$\begin{aligned} \mathbb{E}[Y_i|D_i = 1] - \mathbb{E}[Y_i|D_i = 0] &= \mathbb{E}[Y_i(1)|D_i = 1] - \mathbb{E}[Y_i(0)|D_i = 0] \\ &= \mathbb{E}[Y_i(1)|D_i = 1] - \mathbb{E}[Y_i(0)|D_i = 1] + \mathbb{E}[Y_i(0)|D_i = 1] - \mathbb{E}[Y_i(0)|D_i = 0] \end{aligned}$$

If we have conditional independence:

$$Y_i(0), Y_i(1) \perp D_i | X$$

Then

$$\begin{aligned} \mathbb{E}[Y_i|D_i = 1, X] - \mathbb{E}[Y_i|D_i = 0, X] &= \mathbb{E}[Y_i(1)|D_i = 1, X] - \mathbb{E}[Y_i(0)|D_i = 0, X] \\ &= \mathbb{E}[Y_i(1)|D_i = 1, X] - \mathbb{E}[Y_i(0)|D_i = 1, X] + \mathbb{E}[Y_i(0)|D_i = 1, X] - \mathbb{E}[Y_i(0)|D_i = 0, X] \\ &= \mathbb{E}[Y_i(1)|D_i = 1, X] - \mathbb{E}[Y_i(0)|D_i = 1, X] + \mathbb{E}[Y_i(0)|D_i = 0, X] - \mathbb{E}[Y_i(0)|D_i = 0, X] \\ &= \mathbb{E}[Y_i(1)|D_i = 1, X] - \mathbb{E}[Y_i(0)|D_i = 1, X] \end{aligned}$$

so



$$\mathbb{E}[\mathbb{E}[Y_i|D_i = 1, X] - \mathbb{E}[Y_i|D_i = 0, X]] = \mathbb{E}[\mathbb{E}[Y_i(1)|D_i = 1, X] - \mathbb{E}[Y_i(0)|D_i = 1, X]] = ATT$$

Moreover, from the propensity score theorem, given a propensity score  $p(X)$ , if  $Y_i(0), Y_i(1) \perp D_i|X$  then  $Y_i(0), Y_i(1) \perp D_i|p(X)$ , so we only need to condition in the propensity score for the comparisons.

With the CIA and the propensity score theorem you can estimate  $ATT$  by comparing treated and control units via matching or binning on the propensity score. But, we can also estimate ATEs by weighting by the inverse of the propensity score. Under the CIA,

$$\begin{aligned} \mathbb{E}\left[\frac{Y_i D_i}{p(X)}\right] &= \mathbb{E}\left[\frac{1}{p(X)} D_i Y_i(1)|X\right] \\ &= \mathbb{E}\left[\frac{\mathbb{E}[D_i|X]}{p(X)}\right] \mathbb{E}[Y_i(1)|X] \\ &= \mathbb{E}[\mathbb{E}[Y_i(1)|X]] = \mathbb{E}[Y_i(1)] \end{aligned}$$

A similar calculation shows that  $\mathbb{E}\left[\frac{(1-D_i)Y_i}{1-p(X)}\right] = \mathbb{E}[Y_i(0)]$ . Together, these results imply that we can estimate  $ATT$  with:

$$ATE = \mathbb{E}\left[\frac{Y_i D_i}{p(X_i)} - \frac{Y_i(1-D_i)}{1-p(X_i)}\right]$$

and  $ATT$  as (reweight for the distribution of  $p(X_i)$  on the treated,  $p(X_i)/\Pr(D_i)$ )

$$ATT = \mathbb{E}\left[\frac{Y_i D_i}{\Pr(D_i)} - \frac{Y_i(1-D_i)p(X_i)}{(1-p(X_i))\Pr(D_i)}\right]$$

We are now ready to introduce CS estimator for  $ATT_{g,t}$ :

$$ATT_{g,t} = E\left[\left(\frac{G_g}{E[G_g]} - \frac{\frac{\hat{p}(X)C}{1-\hat{p}(X)}}{E\left[\frac{\hat{p}(X)C}{1-\hat{p}(X)}\right]}\right)(Y_t - Y_{g-1})\right]$$

The first part of the product is introducing the propensity score adjustment, using the generalized propensity score because here we have multiple treatments (the different values of  $g$ ). This is multiplied by the first difference of the outcome instead of the level. So the time difference is unadjusted, and the cross-sectional difference of the time differences is adjusted via the propensity score.

CS's estimators are no different to cross-sectional estimators via the propensity score. The parallel trends assumptions play the role of the CIA assumption in a general setting.

CS also propose outcome-regression and doubly-robust estimators for  $ATT_{g,t}$ . Intuitively, instead of adjusting for the propensity score, we could have adjusted the outcomes via regression. If the regression is correctly specified, under CIA/parallel trends, the estimator is equivalent to the propensity score one. Or, we can both adjust via regression and reweight using the propensity score. This estimator is doubly-robust: we only need either the outcome model or the propensity score model to be correct.

(Discuss aggregation possibilities)

## References

- Angrist, Joshua D and Jörn-Steffen Pischke (2009), *Mostly harmless econometrics: An empiricist's companion*. Princeton university press.
- Borusyak, Kirill and Xavier Jaravel (2018), *Revisiting event study designs*. SSRN.
- Callaway, Brantly and Pedro HC Sant'Anna (2021), "Difference-in-differences with multiple time periods." *Journal of econometrics*, 225, 200–230.
- Goodman-Bacon, Andrew (2021), "Difference-in-differences with variation in treatment timing." *Journal of Econometrics*, 225, 254–277.
- Roth, Jonathan, Pedro HC Sant'Anna, Alyssa Bilinski, and John Poe (2023), "What's trending in difference-in-differences? a synthesis of the recent econometrics literature." *Journal of Econometrics*.