

# Time-Weighted Difference-in-Differences: Accounting for Common Factors in Short $T$ Panels

Timo Schenk\*

February 15, 2022

PRELIMINARY WORK

PLEASE DO NOT CITE OR REDISTRIBUTE WITHOUT PERMISSION

## Abstract

I propose a time-weighted difference-in-differences (TWDID) estimation approach that is robust against time-varying common factors in short  $T$  panels. Time weighting substantially reduces both bias and variance compared to an unweighted DID estimator through balancing the pre-treatment and post-treatment factors. To conduct valid inference on the average treatment effect, I develop a correction term that adjusts conventional standard errors for the presence of weight estimation uncertainty. Revisiting a study on the effect of a cap-and-trade program on NOx emissions, TWDID estimation reduces the standard errors of the estimated treatment effect by 10% compared to a conventional DID approach.

**Keywords:** Synthetic Difference-in-differences, Interactive Fixed Effects, Causal Inference, Panel Data.

## 1 Introduction

The presence of interactive fixed effects, for example due to time-varying common factors, leads to biased difference-in-difference (DID) estimates. While the estimators of [Arkhangelsky, Athey, Hirshberg, Imbens, and Wager \(2021\)](#) and [Chan and Kwok \(2021\)](#) address this issue in large  $T$  panels, the question remains how to account for common factors in short  $T$  panels.

In this paper, I suggest using a time-weighted DID (TWDID) estimator. Factor imbalances between pre-treatment and post-treatment periods cause the DID estimator to be biased. This bias can be eliminated using time

---

\*Department of Economics, University of Amsterdam ([t.d.schenk@uva.nl](mailto:t.d.schenk@uva.nl)). I am grateful for the extensive feedback from my advisors Frank Kleibergen and Andreas Pick as well as helpful discussions during the UvA econometrics seminars.

weights which give a higher weight to pre-treatment periods that are more similar to the post-treatment periods. In theory, a complete elimination of the bias requires oracle weights that perfectly balance the factors. In practice, weights that are estimated from the control unit data succeed in reducing the bias substantially.

A second effect of the factor imbalance is that it amplifies the variance of the DID estimator. By balancing the factors, time weighting reduces the variance and leads to more accurate estimates. In fact, when the number of units is large compared to the number of periods, the estimated weights converge to pseudo-true weights which minimize the variance of the estimated treatment effect. Simulations show that the amount by which the variance is reduced outweighs the additional variance caused by the weight estimation uncertainty.

As a consequence of the bias, inference based on DID estimation will be substantially oversized. The TWDID estimator reduces this problem through bias correction. However, the presence of estimated weights still leads to empirical size in excess of the nominal size when using standard covariance estimators. I propose a two-step standard error procedure that eliminates the remaining size distortions of TWDID inference. First, the weighted cluster-covariance matrix estimator (Arellano, 1987) can be used to estimate the variance of the estimated treatment effect under pseudo-true weights. For the second step, I develop a correction term that accounts for the presence of estimated weights. It uses the fact that in short  $T$  the time weights are asymptotically normal around the pseudo true weights.

Revisiting a study by Deschenes, Greenstone, and Shapiro (2017) on the effect of a cap-and-trade program on NOx emissions, I compare the practical differences between DID and TWDID estimation in a short  $T$  panel. I find evidence of time-varying common factors in the pre-treatment data. While the point estimates only differ to a small degree, TWDID reduces the standard errors of the estimated average treatment effect by 10% compared to a conventional DID approach.

The TWDID estimator is a restricted version of the Arkhangelsky et al. (2021) synthetic DID (SDID) estimator. Next to time weights, the SDID estimator also uses synthetic control unit weights to address the influence of common factors. Consistency results with estimated unit weights require a larger number of pre-treatment periods (Abadie, Diamond, and Hainmueller, 2010; Ferman, 2019; Abadie and L'hour, 2020). Therefore, several challenges occur when applying the SDID estimator in short  $T$  panels. First, the unit weight estimation causes additional variation which might outweigh any additional balancing gains. Second, it is unclear how to conduct inference in presence of estimated unit weights. Arkhangelsky et al. (2021) propose a jackknife based variance estimator which leads to conservative standard errors. However, this result requires that the weight estimation noise is negligible, which only holds for large  $T$ .

Other approaches such as [Chan and Kwok \(2021\)](#) and [Gobillon and Magnac \(2016\)](#) rely on estimating the factor structure by methods of principal components, which also requires large  $T$ . Avoiding the estimation of the whole factor structure, [Pesaran \(2006\)](#) finds cross-sectional averages to sufficiently proxy the factors in each period. In absence of covariates, this approach could account for one common factor when  $T$  is small. When evaluating a one-time policy intervention, however, it is not required to control for the factors in each period. Instead, it suffices to balance the differences between the average pre-treatment and post-treatment factors.

The remainder of the paper is structured as follows. [Section 2](#) covers the Theory. [Section 2.1](#) introduces the interactive fixed effects model and defines the TWDID estimator. [Section 2.2](#) shows the bias and variance reduction properties. [Section 2.3](#) covers inference. [Section 3](#) illustrates the theoretical results with simulations. [Section 4](#) compares TWDID and DID estimation in the study of [Deschenes et al. \(2017\)](#) on NOx emissions.

## 2 Theory

### 2.1 Setting

Consider a policy intervention starting in period  $t = T_0$  which affects units  $i = N_0 + 1, \dots, N$  that are part of a large sample  $i = 1, \dots, N_0, N_0 + 1, \dots, N$  observed over a small number of periods  $t = 1, \dots, T_0, T_0 + 1, \dots, T$ . We seek to estimate the average treatment effect in the post-treatment periods  $t = T_0 + 1, \dots, T$  on the outcome  $y$ . Let  $D_i = \mathbb{I}(i > N_0)$  indicate whether unit  $i$  is ever treated and let  $\mathbb{N}_j = \{i : D_i = j\}$ ,  $j = 0, 1$  be the sets of untreated and treated units, respectively. Let  $\kappa = \frac{N_0}{N}$  denote the share of untreated units and  $\pi = \frac{T_0}{T}$  the share of pre-treatment periods.

We use the interactive fixed model ([Bai, 2009](#))

$$y_{it} = \beta_i + \tau D_{it} + \boldsymbol{\lambda}'_i \mathbf{f}_t + \varepsilon_{it} \quad (1)$$

where  $D_{it} = \mathbb{I}(i > N_0, t > T_0)$  is a dummy that is one for treated units in periods after the treatment started and zero otherwise,  $\tau$  is the homogeneous treatment effect,  $\beta_i$  unit fixed effects,  $\mathbf{f}_t$  and  $\boldsymbol{\lambda}_i$  are  $r$ -dimensional vectors of common factors and loadings, and  $\varepsilon_{it}$  is an idiosyncratic error component.

Unobserved factor structures  $\boldsymbol{\lambda}'_i \mathbf{f}_t$  are present in many economic variables. [Bai \(2009\)](#) discusses, among others, the following examples. In microeconomics,  $\boldsymbol{\lambda}_i$  can be thought of a vector of unobserved, time-invariant characteristics of individual  $i$ . In contrast to the fixed effects  $\beta_i$ , they have a time-varying impact measured by  $\mathbf{f}_t$  on the outcome  $y_{it}$ . In macroeconomics, the factors  $\mathbf{f}_t$  are common shocks (e.g. technology or weather shocks) that have an heterogeneous impact  $\boldsymbol{\lambda}_i$  on unit  $i$ .

I make the following assumptions.

ASSUMPTION 1 (Correlated loadings).  $E[\boldsymbol{\lambda}_i|D_i = 1] - E[\boldsymbol{\lambda}_i|D_i = 0] = \boldsymbol{\xi}_\lambda < \infty$  and  $\text{var}[\boldsymbol{\lambda}_i] = \boldsymbol{\Sigma}_{\lambda,i}$  with  $\lim_{n \rightarrow \infty} \frac{1}{N_j} \sum_{i \in \mathbb{N}_j} \boldsymbol{\Sigma}_{\lambda,i} = \boldsymbol{\Sigma}_\lambda^{(j)}$  for  $j = 0, 1$ , both positive definite  $r \times r$  matrices.

ASSUMPTION 2 (Strict exogeneity). For every  $i$ ,  $E[\boldsymbol{\varepsilon}_i|\beta_i, D_i, \boldsymbol{\lambda}_i, \mathbf{F}] = \mathbf{0}$ . Moreover,  $E[\boldsymbol{\varepsilon}_i \boldsymbol{\varepsilon}_i' | D_i] = \boldsymbol{\Sigma}_{\varepsilon,i}$  with  $\lim_{n \rightarrow \infty} \frac{1}{N_j} \sum_{i \in \mathbb{N}_j} \boldsymbol{\Sigma}_{\varepsilon,i} = \boldsymbol{\Sigma}_\varepsilon^{(j)}$  for  $j = 0, 1$ . Positive definite  $T \times T$  matrices.

ASSUMPTION 3 (Random sampling).  $(\boldsymbol{\varepsilon}_i, \boldsymbol{\lambda}_i)$  are independent over the cross section.  $\kappa, \pi \in (0, 1)$  are both constant as  $N \rightarrow \infty$  and  $T_0 \geq 2$ .

Assumption 1 is the central characteristic of the model. It allows the loadings  $\boldsymbol{\lambda}_i$  to differ systematically between treated and untreated units. The loading imbalance  $\boldsymbol{\xi}_\lambda$  measures how much more the treated units are on average affected by the common factors  $\mathbf{f}_t$ . It also nests the two-way fixed effects model as a special case for  $\boldsymbol{\xi}_\lambda = 0$  and  $\boldsymbol{\Sigma}_\lambda^{(1)} = \boldsymbol{\Sigma}_\lambda^{(0)} = 0$ , since then  $\boldsymbol{\lambda}_i = \boldsymbol{\lambda}$  for all  $i$ . In that case the factor structure  $\boldsymbol{\lambda}_i' \mathbf{f}_t$  reduces to a time fixed effect  $\gamma_t = \boldsymbol{\lambda}' \mathbf{f}_t$ .

Under Assumption 2, the treatment assignment is strictly exogenous once conditioned on the loadings, fixed effects and the factors. Moreover, I allow for heteroskedasticity and arbitrary serial dependence of the idiosyncratic errors. The common factors  $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_T)'$  ( $T \times r$ ) are realizations of an arbitrary deterministic or stochastic process. The number of factors  $r$  is unknown and fixed. Without loss of generality I assume that  $\mathbf{F}' \iota = \mathbf{0}$ . All results hold conditional on  $\mathbf{F}$ .

Assumption 3 imposes independence of the error component over the cross section. It requires the number of treated and untreated units to grow at the same rate. We need at least two pre-treatment periods.

I propose a restricted version of the [Arkhangelsky et al. \(2021\)](#) synthetic DID estimator which uses only time weights and refer to it as the time-weighted DID estimator. It is computed in two steps.

1. Obtain a  $T_0$  vector of time weights  $\hat{\mathbf{v}} = (\hat{v}_1, \dots, \hat{v}_{T_0})'$  using the de-meaned outcomes  $\hat{y}_{i,t} = y_{i,t} - \frac{1}{N_0} \sum_{i=1}^{N_0} y_{i,t}$  of the untreated units only. Regress the average post-treatment outcome  $\bar{y}_{i,[1]}$  on the pre-treatment outcomes  $\dot{\mathbf{y}}_{i,[0]} = (\dot{y}_{i,1}, \dots, \dot{y}_{i,T_0})'$

$$\hat{\mathbf{v}} = \arg \min_{\mathbf{v} \in \mathbb{V}} \sum_{i=1}^{N_0} (\bar{y}_{i,[1]} - \dot{\mathbf{y}}_{i,[0]}' \mathbf{v})^2 \quad (2)$$

with  $\mathbb{V} = \{\mathbf{v} \in \mathbb{R}^{T_0} : v_t \geq 0, \sum_{t=1}^{T_0} v_t = 1\}$  the set of non-negative weights that sum to one.

2. Obtain the time-weighted DID estimator  $\hat{\tau}(\hat{\mathbf{v}})$  as solution to the weighted two-way fixed effect regression

$$\min_{\tau, \mu, \gamma} \sum_{i=1}^N \sum_{t=1}^T v_t (y_{it} - \tau D_{it} - \mu_i - \gamma_t)^2 \quad (3)$$

The resulting estimator is

$$\hat{\tau}(\mathbf{v}) = \bar{\Delta}_{[1]} - \sum_{t=1}^{T_0} v_t \Delta_t \quad (4)$$

with  $\Delta_t = \bar{y}_t^{(1)} - \bar{y}_t^{(0)}$ ,  $\bar{\Delta}_{[1]} = \frac{1}{T-T_0} \sum_{t>T_0} \Delta_t$ ,  $\bar{y}_t^{(1)} = \frac{1}{N_1} \sum_{i=N_0+1}^N y_{it}$  and  $\bar{y}_t^{(0)} = \frac{1}{N_0} \sum_{i=1}^{N_0} y_{it}$  the treated and untreated units' average outcome, respectively.

The DID estimator is the special case of the TWDID estimator with equal weights  $\bar{\mathbf{v}} = \frac{\iota_{T_0}}{T_0}$ . The synthetic DID estimator of [Arkhangelsky et al. \(2021\)](#) uses both unit weights and time weights in product form  $\omega_i v_t$  in (3). The treated units are weighted equally and the weights of the untreated units are non-negative and sum to one. The estimator has the form of  $\hat{\tau}(\mathbf{v})$  with simple averages over the control units  $\bar{y}_t^{(0)}$  replaced by weighted averages  $\sum_{i \in \mathbb{N}_0} \omega_i y_{it}$ . The control unit weights are estimated from the pre-treatment data in the same manner as the pre-treatment time weights are estimated from the control unit data. Unit weight estimation, however, requires an additional penalty term to ensure uniqueness when  $N_0$  is large compared to  $T_0$ .

## 2.2 Bias and variance reduction through factor balancing

In the following part I will first consider the properties of  $\hat{\tau}(\mathbf{v})$  under fixed weights to document the issue of factor imbalance. I then show how estimated time weights reduce the bias and the asymptotic variance of the estimated treatment effect compared to the unweighted DID approach.

Consider the treatment effect estimate  $\hat{\tau}(\mathbf{v})$  for a given vector of weights  $\mathbf{v}$ . For equal weights  $\mathbf{v} = \bar{\mathbf{v}}$  this resembles the DID estimator.

**THEOREM 1.** *Suppose Assumptions 1-3 hold. Then for every  $\mathbf{v} \in \mathbb{V}$ ,*

1.  $E[\hat{\tau}(\mathbf{v}) | \mathbf{F}] = \tau + b(\mathbf{v})$  with bias  $b(\mathbf{v}) = \boldsymbol{\xi}'_{\lambda} \boldsymbol{\xi}_f(\mathbf{v})$  and the weighted factor imbalance  $\boldsymbol{\xi}_f(\mathbf{v}) = \mathbf{f}_{(1)} - \mathbf{F}'\mathbf{v}$ ,
2.  $N \text{var}[\hat{\tau}(\mathbf{v}) | \mathbf{F}] = \boldsymbol{\xi}_f(\mathbf{v})' \boldsymbol{\Sigma}_{\lambda} \boldsymbol{\xi}_f(\mathbf{v}) + V_{\varepsilon}(\mathbf{v}) =: V_{\hat{\tau}}(\mathbf{v})$  with  $\boldsymbol{\Sigma}_{\lambda} = \text{var}[\bar{\boldsymbol{\lambda}}^{(1)} - \bar{\boldsymbol{\lambda}}^{(0)}]$ ,  $V_{\varepsilon}(\mathbf{v}) = \text{var}[\bar{\Delta}_{\varepsilon, [1]} - \sum_{t \leq T_0} v_t \Delta_{\varepsilon, t}]$ ,  $\Delta_{\varepsilon, t} = \bar{\varepsilon}_t^{(1)} - \bar{\varepsilon}_t^{(0)}$ , and
3.  $\sqrt{N}(\hat{\tau}(\mathbf{v}) - \tau - b(\mathbf{v})) \xrightarrow{d} N[0, V_{\hat{\tau}}(\mathbf{v})]$  as  $N \rightarrow \infty$ .

The weighted factor imbalance  $\boldsymbol{\xi}_f(\mathbf{v})$  will play an important role. It is the difference between the average post-treatment factors and the weighted average pre-treatment factors and affects the estimated treatment effect in two ways. First, the combination of a non-zero loading imbalance  $\boldsymbol{\xi}_\lambda$  and factor imbalance  $\boldsymbol{\xi}_f(\mathbf{v})$  leads to a first order bias term  $b(\mathbf{v})$ . For example, consider the case with one common factor  $f_t$ , which affects treated units are on average more than untreated units ( $\xi_\lambda > 0$ ). If  $f_t$  is on average in the post-treatment periods higher than in the pre-treatment periods,  $\hat{\tau}(\mathbf{v})$  will overestimate the treatment effect. Second, it increases the part of the variance resulting from variation in the loadings. This holds irrespective of whether the treatment assignment  $D_i$  correlates with the loadings  $\lambda_i$ , as long as they have within group variation  $\boldsymbol{\Sigma}_\lambda > 0$ .

The properties of the DID estimator follow as the special case of equal weights  $\mathbf{v} = \bar{\mathbf{v}}$ . Researches typically refers to the common trend assumption as condition for unbiasedness. In the current setting, a trend means a non-zero factor imbalance  $\boldsymbol{\xi}_f(\bar{\mathbf{v}})$ . The trends are common if the factors affect treated and untreated units equally. Hence the DID estimator is unbiased ( $b(\bar{\mathbf{v}}) = 0$ ) if either the trends are common ( $\boldsymbol{\xi}_\lambda = 0$ ) or there are no trends ( $\boldsymbol{\xi}_f(\bar{\mathbf{v}}) = 0$ ).

Now consider weights  $\hat{\mathbf{v}}$  estimated from the control unit data as per (2). As the number of control units  $N_0$  grows, they converge in probability to the pseudo-true time weights  $\mathbf{v}^*$  which solve the population equivalent of (2)

$$\mathbf{v}^* = \arg \min_{\mathbf{v} \in \mathbb{V}} \left\{ \boldsymbol{\xi}_f(\mathbf{v})' \boldsymbol{\Sigma}_\lambda^{(0)} \boldsymbol{\xi}_f(\mathbf{v}) + V_\varepsilon^{(0)}(\mathbf{v}) \right\} \quad (5)$$

with  $V_\varepsilon^{(0)}(\mathbf{v}) = \text{var}[\bar{\varepsilon}_{[1]}^{(0)} - \sum_{t \leq T_0} v_t \bar{\varepsilon}_t^{(0)}]$ . The pseudo-true weights minimize an expression close to the variance of  $\hat{\tau}(\mathbf{v})$  derived in Theorem 1. It is not just influenced by the factor imbalance  $\boldsymbol{\xi}_f(\mathbf{v})$ , but also by the error variance  $V_\varepsilon^{(0)}(\mathbf{v})$ . As a consequence, the pseudo-true weights do generally not balance the factors entirely. The following Theorem establishes asymptotic normality around  $\mathbf{v}^*$ .

**THEOREM 2.** *Suppose Assumptions 1-3 hold. Let  $T_+$  be the number of strictly positive elements of  $\mathbf{v}^*$  with  $1 \leq T_+ \leq T_0$ . As  $N_0 \rightarrow \infty$ ,*

1.  $\hat{\mathbf{v}} \xrightarrow{p} \mathbf{v}^*$
2.  $\sqrt{N}(\hat{\mathbf{v}} - \mathbf{v}^*) \xrightarrow{d} N[\mathbf{0}, \frac{1}{\kappa} \boldsymbol{\Sigma}_{\hat{\mathbf{v}}}]$  with  $\text{rk} \boldsymbol{\Sigma}_{\hat{\mathbf{v}}} = T_+ - 1$ ,  $(\boldsymbol{\Sigma}_{\hat{\mathbf{v}}})_{s,t} = 0$  if  $\min\{v_s^*, v_t^*\} = 0$ .

Appendix A.2 contains the proof.

Most importantly, the estimated weights  $\hat{\mathbf{v}}$  and the treatment effect estimate  $\hat{\tau}(\mathbf{v})$  converge at the same rate  $\frac{1}{\sqrt{N}}$ . Using estimated weights, the estimation error then becomes

$$\hat{\tau}(\hat{\mathbf{v}}) - \tau - b(\mathbf{v}^*) = Z_N(\mathbf{v}^*) - (\mathbf{F}_{[0]} \boldsymbol{\xi}_\lambda)'(\hat{\mathbf{v}} - \mathbf{v}^*) + O_p\left(\frac{1}{N}\right)$$

with  $Z_N(\mathbf{v}^*) = \hat{\tau}(\mathbf{v}^*) - \tau - b(\mathbf{v}^*)$  the estimation error under pseudo-true weight and  $\mathbf{F}_{[0]} = (\mathbf{f}_1, \dots, \mathbf{f}_{T_0})'$  the  $(T_0 \times r)$  matrix of pre-treatment factors. This leads us to the following result

**THEOREM 3.** *Suppose Assumptions 1-3 hold. Then, as  $N \rightarrow \infty$ ,*

1.  $\hat{\tau}(\hat{\mathbf{v}}) \xrightarrow{p} \tau + b(\mathbf{v}^*)$
2.  $\sqrt{N}(\hat{\tau}(\hat{\mathbf{v}}) - \tau - b(\mathbf{v}^*)) \xrightarrow{d} N[0, V_{\hat{\tau}(\hat{\mathbf{v}})}]$

with

$$V_{\hat{\tau}(\hat{\mathbf{v}})} = V_{\hat{\tau}(\mathbf{v}^*)} + \frac{1}{\kappa} (\mathbf{F}_{[0]} \boldsymbol{\xi}_\lambda)' \boldsymbol{\Sigma}_{\hat{\mathbf{v}}} \mathbf{F}_{[0]} \boldsymbol{\xi}_\lambda$$

The magnitude of the bias  $b(\mathbf{v}^*) = \boldsymbol{\xi}_\lambda' \boldsymbol{\xi}_f(\mathbf{v}^*)$  depends on the remaining factor imbalance under pseudo-true weights. The limit variance consists of two parts. First,  $V_{\hat{\tau}(\mathbf{v}^*)}$  is the variance of the treatment effect estimator under fixed, pseudo-true weights  $\mathbf{v}^*$ . The second term comes from the weight estimation noise and is only present if  $\boldsymbol{\xi}_\lambda \neq 0$ . In this case  $\hat{\tau}(\hat{\mathbf{v}})$  does not have oracle properties as its asymptotic distribution differs from the one of  $\hat{\tau}(\mathbf{v}^*)$ .

Comparing the DID estimator  $\hat{\tau}(\bar{\mathbf{v}})$  to the time-weighted version  $\hat{\tau}(\hat{\mathbf{v}})$  leads to the following conclusions. The bias of the latter is smaller if the pseudo-true weights decrease the factor imbalance compared to equal weights. This is arguably the case in relevant scenarios, although technically it is possible to make up cases in which it does not hold. Next, weighting has a two-fold effect on the relative variance

$$\frac{V_{\hat{\tau}(\hat{\mathbf{v}})}}{V_{\hat{\tau}(\bar{\mathbf{v}})}} = \frac{V_{\hat{\tau}(\mathbf{v}^*)}}{V_{\hat{\tau}(\bar{\mathbf{v}})}} + \frac{\mathbf{F}_{[0]} \boldsymbol{\xi}_\lambda' \boldsymbol{\Sigma}_{\hat{\mathbf{v}}} \mathbf{F}_{[0]} \boldsymbol{\xi}_\lambda}{\kappa V_{\hat{\tau}(\bar{\mathbf{v}})}}$$

First,  $\mathbf{v}^*$  minimizes the first term by construction. This comes at the cost of weight estimation noise, hence the second term. The following corollary provides high-level sufficient conditions to assure that TWDID in fact reduces bias and variance.

**COROLLARY 1.** *Suppose Assumptions 1-3 hold. If  $\boldsymbol{\Sigma}_\lambda^{(j)}, \boldsymbol{\Sigma}_\varepsilon^{(j)}$ ,  $j = 0, 1$  and  $\mathbf{F}$  are such that*

1.  $\boldsymbol{\xi}_\lambda' (\boldsymbol{\xi}_f(\bar{\mathbf{v}}) \boldsymbol{\xi}_f(\bar{\mathbf{v}})' - \boldsymbol{\xi}_f(\mathbf{v}^*) \boldsymbol{\xi}_f(\mathbf{v}^*)') \boldsymbol{\xi}_\lambda > 0$ ,
2.  $\frac{V_{\hat{\tau}(\mathbf{v}^*)}}{V_{\hat{\tau}(\bar{\mathbf{v}})}} < 1 - \frac{\mathbf{F}_{[0]} \boldsymbol{\xi}_\lambda' \boldsymbol{\Sigma}_{\hat{\mathbf{v}}} \mathbf{F}_{[0]} \boldsymbol{\xi}_\lambda}{\kappa V_{\hat{\tau}(\bar{\mathbf{v}})}}$

*then the TWDID estimator  $\hat{\tau}(\hat{\mathbf{v}})$  has a smaller bias and a smaller variance than the DID estimator  $\hat{\tau}(\bar{\mathbf{v}})$ .*

The Monte Carlo simulations in Section 3 show that TWDID substantially reduces bias and variance in a setting with one gaussian factor.

### 2.3 Inference with two-step standard errors

I propose a consistent estimator of the limit variance derived in Theorem 3

$$\widehat{V}_{\hat{\tau}(\hat{\mathbf{v}})} = \widehat{V}_{\hat{\tau}}(\hat{\mathbf{v}}) + \frac{1}{\kappa} \dot{\Delta}'_{[0]} \widehat{\Sigma}_{\hat{\mathbf{v}}} \dot{\Delta}_{[0]} \quad (6)$$

which consists of two parts. The first part  $\widehat{V}_{\hat{\tau}}(\hat{\mathbf{v}})$  is the cluster-covariance robust variance estimator of Arellano (1987) applied to the weighted data. It consistently estimates the variance under pseudo-true weights  $V_{\hat{\tau}}(\mathbf{v}^*)$ . The second part accounts for the additional variance caused by the weight estimation noise. It consist of the demeaned average pre-treatment differences  $\dot{\Delta}_{[0]} = (\Delta_1 - \bar{\Delta}_{[0]}, \dots, \Delta_{T_0} - \bar{\Delta}_{[0]})'$  with  $\bar{\Delta}_{[0]} = \frac{1}{T_0} \sum_{t \leq T_0} \Delta_t$  and a consistent estimator  $\widehat{\Sigma}_{\hat{\mathbf{v}}}$  of the weight variance. In the remainder of this section I will explain how to construct the different components of  $\widehat{V}_{\hat{\tau}(\hat{\mathbf{v}})}$  and associated confidence intervals in more detail.

The first part of the estimated variance,  $\widehat{V}_{\hat{\tau}}(\hat{\mathbf{v}})$ , is obtained in the following way. First estimate the time weights  $\hat{\mathbf{v}}$ . Next, weigh only the pre-treatment outcomes and call them  $\tilde{y}_{it} = T_0 \hat{v}_t y_{it}$  for  $t \leq T_0$  and  $\tilde{y}_{it} = y_{it}$  for  $t > T_0$ . Run a two-way fixed effects regression of the weighted outcomes  $\tilde{y}_{it}$  on the treatment indicator  $D_{it}$ , which yields  $\hat{\tau}(\hat{\mathbf{v}})$ . Applying the Arellano (1987) cluster-covariance estimator<sup>1</sup> on the weighted data then provides  $\widehat{V}_{\hat{\tau}}(\hat{\mathbf{v}})$ . The following Theorem ensures consistency. Its details are explained in the Appendix.

**THEOREM 4.** *Suppose Assumptions 1-3 hold. Then for any  $\mathbf{v} \in \mathbb{V}$*

$$\widehat{V}_{\hat{\tau}}(\mathbf{v}) \xrightarrow{p} V_{\hat{\tau}}(\mathbf{v})$$

as  $N \rightarrow \infty$ . Moreover,

$$\widehat{V}_{\hat{\tau}}(\hat{\mathbf{v}}) \xrightarrow{p} V_{\hat{\tau}}(\mathbf{v}^*)$$

with  $\hat{\mathbf{v}}$  the estimated weights as per (2) and  $\mathbf{v}^*$  their probability limit as defined in (5).

Next, I will explain why  $\dot{\Delta}'_{[0]} \widehat{\Sigma}_{\hat{\mathbf{v}}} \dot{\Delta}_{[0]}$  provides an adequate estimator of  $(\mathbf{F}_{[0]} \boldsymbol{\xi}_{\lambda})' \Sigma_{\hat{\mathbf{v}}} \mathbf{F}_{[0]} \boldsymbol{\xi}_{\lambda}$ . First, the demeaned average pre-treatment differences become

$$\dot{\Delta}_{[0]} = \mathbf{F}_{[0]} \boldsymbol{\xi}_{\lambda} + O_p\left(\frac{1}{\sqrt{N}}\right)$$

implying that they consistently estimate  $\mathbf{F}_{[0]} \boldsymbol{\xi}_{\lambda}$  as  $N \rightarrow \infty$ . The estimator of the weight variance  $\widehat{\Sigma}_{\hat{\mathbf{v}}}$  can be explained as follows. Let  $\hat{\mathbf{v}}$  be the  $T_0$  vector of estimated weights as per (2). As Theorem 3 establishes, only the non-zero weights matter for the limiting weight variance  $\Sigma_{\hat{\mathbf{v}}}$ . We can write

<sup>1</sup>Often referred to as clustering the standard errors on the unit (or individual) level, for instance in Bertrand, Duflo, and Mullainathan (2004)



the non-zero weights as an unrestricted least-squares estimate and thus use least-squares type of standard errors to estimate its variance.

Let  $\hat{\boldsymbol{v}}_{[+]}$  be the  $T_+$  vector which just contains the positive weights. Because the weights sum to one, I can write

$$\hat{\boldsymbol{v}}_{[+]} = \boldsymbol{e}_1 + \boldsymbol{R}\hat{\boldsymbol{v}}_{-1}$$

with  $\boldsymbol{e}_1$  the  $T_+$ -dimensional unit vector,  $\boldsymbol{R} = \begin{pmatrix} -\boldsymbol{1}'_{T_+-1} \\ \boldsymbol{I}_{T_+-1} \end{pmatrix}$  a  $T_+ \times T_+ - 1$  matrix and  $\hat{\boldsymbol{v}}_{-1}$  the last  $T_+ - 1$  elements of  $\hat{\boldsymbol{v}}_{[+]}$ . The latter can be written as the unrestricted least-squares estimate

$$\hat{\boldsymbol{v}}_{-1} = (\tilde{\boldsymbol{Y}}'_{[0]}\tilde{\boldsymbol{Y}}_{[0]})^{-1}\tilde{\boldsymbol{Y}}'_{[0]}\tilde{\boldsymbol{y}}_{[1]}$$

with  $\tilde{\boldsymbol{Y}}_{[0]} = \dot{\boldsymbol{Y}}_{[+]}^{(0)}\boldsymbol{R}$ ,  $\tilde{\boldsymbol{y}}_{[1]} = \dot{\boldsymbol{y}}_{[1]}^{(0)} - \dot{\boldsymbol{Y}}_{[+]}^{(0)}\boldsymbol{e}_1$  and  $\dot{\boldsymbol{Y}}_{[+]}^{(0)}$  the  $N_0 \times T_+$  matrix of demeaned outcomes of the control units in the remaining pre-treatment periods. Its variance can be estimated by

$$\hat{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{v}}_{-1}} = Q_N(\hat{\boldsymbol{v}})(\tilde{\boldsymbol{Y}}'_{[0]}\tilde{\boldsymbol{Y}}_{[0]})^{-1}$$

with  $Q_N(\hat{\boldsymbol{v}})$  the mean sum of squared residuals of the time weight estimation. Finally, the estimator of the weight covariance matrix is

$$\hat{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{v}}_{[+]}} = Q_N(\hat{\boldsymbol{v}})\boldsymbol{R}(\tilde{\boldsymbol{Y}}'_{[0]}\tilde{\boldsymbol{Y}}_{[0]})^{-1}\boldsymbol{R}' \quad (7)$$

which follows from  $\text{var}[\hat{\boldsymbol{v}}_{[+]}] = \boldsymbol{R}\text{var}[\hat{\boldsymbol{v}}_{-1}]\boldsymbol{R}'$ .

Inference can be based on the t-statistic following

$$\text{P} \left( \left| \frac{\hat{\tau}(\hat{\boldsymbol{v}}) - \tau_0 - b(\boldsymbol{v}^*)}{\sqrt{\hat{V}_{\hat{\tau}(\hat{\boldsymbol{v}})}/N}} \right| > q_{1-\alpha/2} \right) \xrightarrow{p} \alpha$$

with  $\tau_0$  the treatment effect under the null,  $\alpha$  the intended size and  $q_{1-\alpha/2}$  the  $1 - \alpha/2$  quantile of the standard normal distribution. Without any further restrictions, the resulting confidence interval

$$\left[ \hat{\tau}(\hat{\boldsymbol{v}}) \pm q_{1-\alpha/2} \sqrt{\hat{V}_{\hat{\tau}(\hat{\boldsymbol{v}})}/N} \right]$$

will be centered around  $\tau_0 + b(\boldsymbol{v}^*)$ . Yet it is more reliable and shorter compared to an unweighted DID approach.

### 3 Monte Carlo Experiments

#### 3.1 Implementation

In each Monte Carlo replication  $r = 1, \dots, R$  I generate data from

$$y_{it}^{(r)} = \tau D_{it} + \lambda_i^{(r)} f_t^{(r)} + \varepsilon_{it}^{(r)}$$

with  $f_t^{(r)} \sim N[0, \sigma_f^2]$  and  $\varepsilon_{it}^{(r)} \sim N[0, 1]$  for  $i = 1, \dots, N$  and  $t = 1, \dots, T$ , all mutually independent. The loadings are  $\lambda_i^{(r)} = \frac{\xi_\lambda}{\sqrt{N}} D_i + \nu_i^{(r)}$  with  $\nu_i^{(r)} \sim N[0, 1]$ . The true treatment effect is  $\tau = 0$ . The number of units is  $N = 100$  of which half are treated ( $\kappa = \frac{N_0}{N} = 0.5$ ) and there is one treated period  $T = T_0 + 1$ . I vary the number of pre-treatment periods  $T_0 \in \{6, 30\}$  and the loading imbalance  $\xi_\lambda \in \{2, 5\}$  along a grid of factor standard deviations  $\sigma_f \in \{0, 0.1, \dots, 2\}$ . For each combination of parameters and sample size I conduct  $R = 10,000$  replications. In each replication I compute the pseudo-true time weights  $\mathbf{v}^*$  as of (5), the estimated time weights as of (2), the DID estimator  $\hat{\tau}_{\text{did}} = \hat{\tau}(\bar{\mathbf{v}})$ , the TWDID estimator  $\hat{\tau}_{\text{twdid}} = \hat{\tau}(\hat{\mathbf{v}})$  and the corresponding variance estimators  $\hat{V}_{\text{did}}, \hat{V}_{\text{twdid}}$ . I compare the following Monte Carlo statistics to assess the performance of the point estimators and inference.

1. I measure the magnitude of the bias term  $b(\mathbf{v}) = \xi_\lambda \xi_f(\mathbf{v})$  in terms of its standard deviation  $\text{sd}[b(\mathbf{v})]$ , with  $\mathbf{v} = \bar{\mathbf{v}}$  for the DID and  $\mathbf{v} = \mathbf{v}^*$  for TWDID.
2. Then I look at the simulated conditional standard deviation  $\text{sd}[\hat{\tau}(\mathbf{v})|\mathbf{F}]$  of the point estimates. I compute it as the Monte Carlo standard deviation of the bias corrected estimator  $\hat{\tau}(\mathbf{v}) - b(\mathbf{v})$ .
3. Finally, I consider the feasible t-statistics of both estimation approaches

$$t_{\text{did}} = \frac{\hat{\tau}_{\text{did}} - \tau_0}{\sqrt{\hat{V}_{\text{did}}/N}}, \quad t_{\text{twdid}} = \frac{\hat{\tau}_{\text{twdid}} - \tau_0}{\sqrt{\hat{V}_{\text{twdid}}/N}}$$

I compute the true rejection frequency under  $\tau_0 = 0$  using critical values from the standard normal distribution.

### 3.2 Results

The simulation results confirm the bias and variance reduction properties of TWDID compared to DID. The top panel of Figure 1 plots the magnitude of the bias against the strength of the factors. The bias of DID increases as the factors become stronger. The bias of TWDID is about 50% lower. This effect is stronger in a setting with a higher loading imbalance (left vs. right). The bottom panel shows the conditional standard deviation, which is independent of the loading imbalance and instead driven by the loading variation. For sufficiently strong factors, the TWDID estimator has a substantially lower standard deviation. For weak factors, however, there is not much to be gained and the weight estimation noise leads to a slightly larger standard deviation.

Consider now inference based on the t-statistics  $t_{\text{did}}$  and  $t_{\text{twdid}}$ . Figure 2 shows the rejection frequency of the standard t-test under the null of no

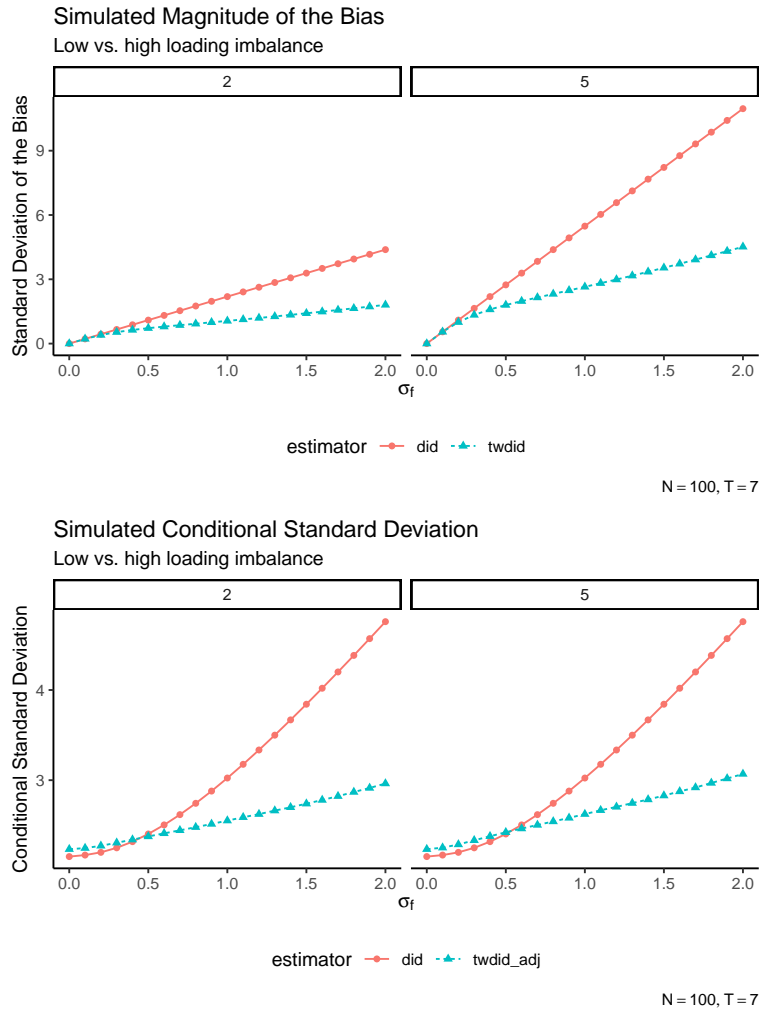


Figure 1: Magnitude of the conditional bias term  $b(\mathbf{v}^*) = \xi_\lambda \xi_f(\mathbf{v}^*)$  (top) and simulated conditional standard deviation  $\text{sd}[\hat{\tau}(\mathbf{v})]$  (bottom) of the DID and TWDID estimator, depending on the factor standard deviation  $\sigma_f$ .

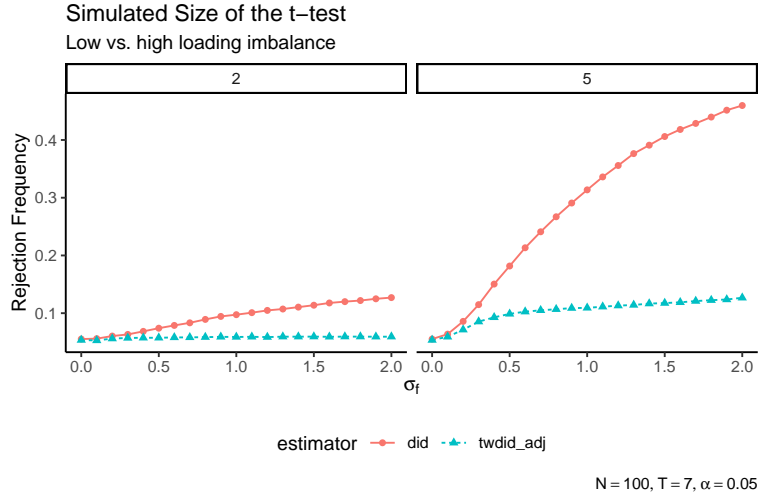


Figure 2: Simulated rejection frequency based on the feasible t-statistic  $t_{\text{did}}$  and  $t_{\text{twdid}}$  depending on the factor strength  $\sigma_f$ . The standard errors of the TWDID estimator have been adjusted for weight estimation uncertainty.

treatment effect  $H_0: \tau = 0$ . If the loading imbalance is sufficiently small (left), the remaining bias of the TWDID estimator is negligible as the test remains size-correct when the factors become stronger. The larger bias of the DID estimator leads its test statistic to be distorted. For large loading imbalances (right), the bias becomes so large that both tests are distorted. Yet this distortion is less severe for TWDID than it is for DID.

## 4 The effect of the NOx Budget Trading Program revisited

I revisit [Deschenes et al. \(2017\)](#) studying the effect of the NOx Budget Trading Program (NBP) 2003-2008 on NOx emissions. It entailed a cap and trade program to reduce NOx emissions from power plants. It was only active in the summer months May - September in the years 2003-2008 in 19 states in the US. In 2003 the program was active only in a subset of the 19 treated states. States not adjacent to the NBP states remain as non-treated states (22 in total).

Data on NOx emissions is available on county level for  $N = 2539$  counties from 1997-2007. We observe  $N_1 = 1,354$  counties in the treated states and  $N_0 = 1,185$  in the untreated states. Per county and year we observe data for the seasons summer and winter, where summer is defined as May - September.

## 4.1 Econometric Specification

Consider the interactive fixed effect model

$$y_{ist} = \tau D_{ist} + \mu_{it} + \nu_{is} + \boldsymbol{\lambda}'_i \mathbf{f}_{st} + \tilde{\varepsilon}_{ist}$$

with  $D_{ist} = \mathbf{I}(i \in \mathcal{N}_1, t > T_0, s = 1)$  an indicator whether NBP is operating in county  $i$  in season  $s = 0, 1$  (winter, summer) of year  $t$ .  $\mu_{it}, \nu_{is}$  are county-year and county-season fixed effects, respectively.  $\mathbf{f}_{st}$  are season-year specific common shocks that affect the emissions of county  $i$  with intensity  $\boldsymbol{\lambda}_i$ .  $\tilde{\varepsilon}_{ist}$  is an idiosyncratic error term. The special case  $\boldsymbol{\lambda}_i = \boldsymbol{\lambda}$  resembles the additive fixed effect model that [Deschenes et al. \(2017\)](#) assume. In that case the factor structure reduces to a season-year fixed effect.

To identify  $\tau$ , eliminate  $\mu_{it}$  by taking the difference between summer and winter observations

$$\check{y}_{it} := y_{i1t} - y_{i0t} = \tau D_{it} + \beta_i + \boldsymbol{\lambda}'_i \check{\mathbf{f}}_t + \varepsilon_{it}$$

with  $\beta_i = \nu_{i1} - \nu_{i0}$ ,  $\check{\mathbf{f}}_t = \mathbf{f}_{1t} - \mathbf{f}_{0t}$  and  $\varepsilon_{it} = \tilde{\varepsilon}_{i1t} - \tilde{\varepsilon}_{i0t}$ . A key assumption hidden in this specification is that the program does not affect emissions in the winter months in the treated years. The identifying assumption is that there exists a vector of weights  $\mathbf{v}_0$  such that

$$\frac{1}{T_1} \sum_{t>T_0} \check{\mathbf{f}}_t = \sum_{t \leq T_0} v_{0,t} \check{\mathbf{f}}_t$$

That is, the average post-treatment factor can be written as a weighted average of pre-treatment factors.

## 4.2 Evidence for common factors

I first obtain evidence against  $\boldsymbol{\lambda}_i = \boldsymbol{\lambda}$  by considering how the difference in average NOx emissions  $\Delta_t = \bar{y}_t^{(1)} - \bar{y}_t^{(0)}$  has evolved prior to the intervention. We can write

$$\Delta_t = \bar{\beta}^{(1)} - \bar{\beta}^{(0)} + \boldsymbol{\xi}'_{\lambda} \mathbf{f}_t + \tau \mathbf{I}(t > T_0) + O_p\left(\frac{1}{\sqrt{N}}\right)$$

where  $\boldsymbol{\xi}_{\lambda} = 0$  under the equal loading assumption. Then, for large  $N$ ,  $\Delta_t$  should be constant prior to the treatment. However, [Figure 3](#) does show variation of  $\Delta_t$  in periods  $t \leq T_0$ .

## 4.3 Estimation Results

I estimate  $\tau$  with a time weighted DID regression as defined in [\(4\)](#). Let  $\hat{\tau}(\mathbf{v})$  be the point estimate for a given  $T_0$  vector of time weights  $\mathbf{v} \in \mathbb{V}$ . I compare three sets of weights  $\mathbf{v}$ . First, I use equal weights  $\bar{\mathbf{v}}$ , which return the DID

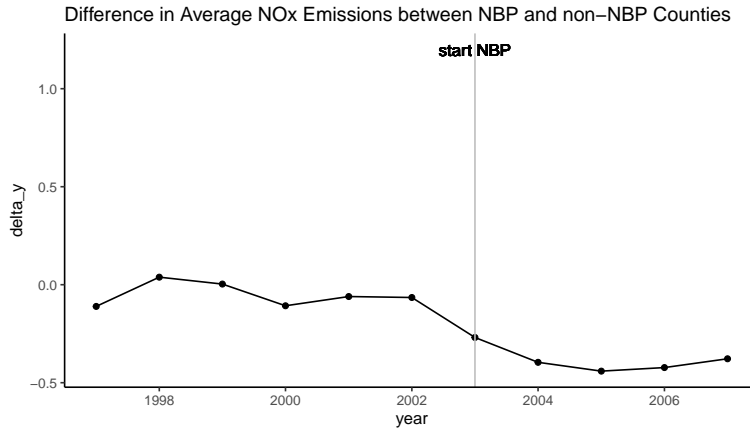


Figure 3: Difference in average NOx emissions  $\bar{y}_t^{(1)} - \bar{y}_t^{(0)}$

estimator used by [Deschenes et al. \(2017\)](#). Second, for the TWDID estimator I use estimated weights  $\hat{v}$  obtained from the restricted regression (2) of pre-treatment emissions on the average post-treatment emission. Third, I consider the TWDID estimator using unrestricted weights  $\hat{v}_u$ , where the only difference to  $\hat{v}$  is that I allow the weights to be non-negative. I obtain the standard errors For the DID estimator I use the cluster-covariance matrix estimator. I obtain the standard errors of the TWDID estimators with the two-step procedure discussed in Section 2.3, which accounts for the presence of estimated weights. I omit the year 2003 from the estimation to circumvent issues related to the staggered implementation.

The estimated time weights are plotted in Figure 4. The last two pre-treatment periods receive almost all the weight, indicating that the common factors of the NOx emissions were in these two years closest to their post-treatment average. The bottom panel of Figure 4 shows the resulting treatment effect estimates and their 95% confidence intervals. All estimators suggest a significant negative effect of the NBP program on NOx emissions. Time weighting leads, in absolute terms, to a slightly lower point estimate (DID: -0.36 vs TWDID: -0.34). The standard errors of both TWDID estimates (0.050) are about 10% lower compared to DID (0.056), and the resulting confidence intervals are more narrow. This results is in line with the variance reduction property of TWDID estimation.

## 5 Conclusion

This paper proposes a time-weighted difference-in-difference (TWDID) estimation approach that is robust against interactive fixed effects in short  $T$  panels. It covers settings with few time periods, many cross-sectional units and sharp treatment timing. Using time weights estimated from the

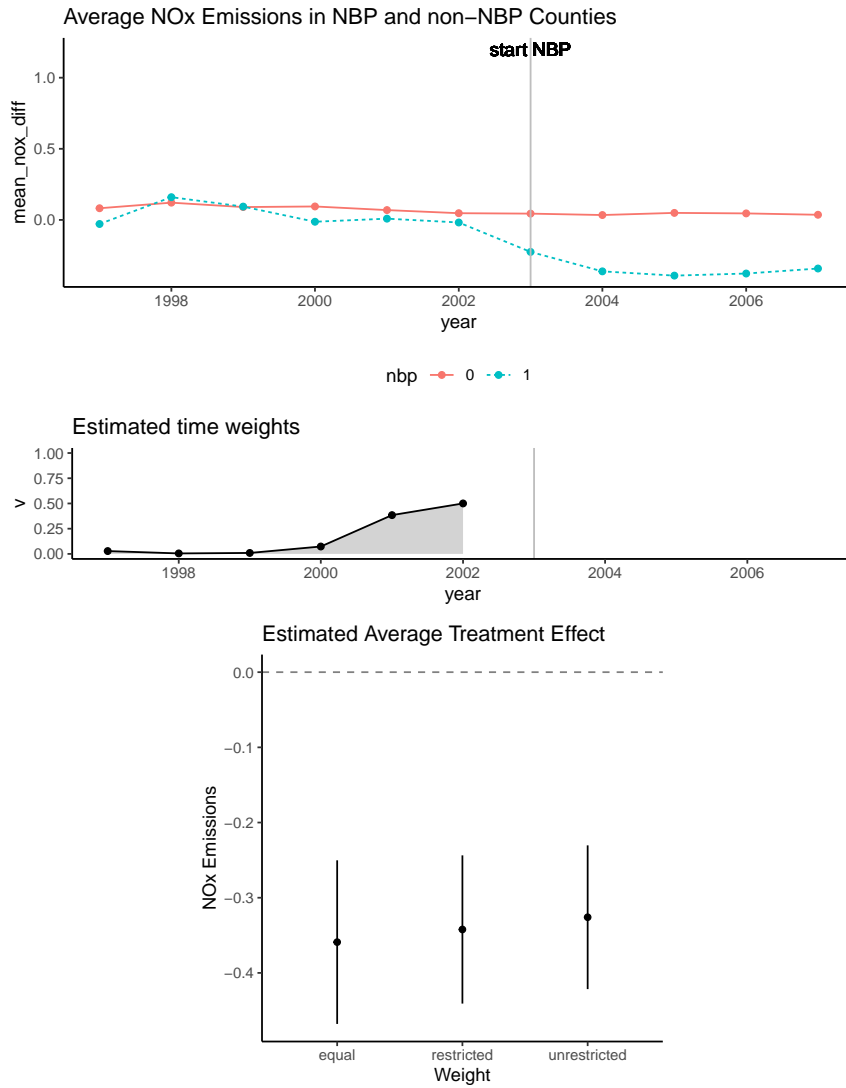


Figure 4: Top: Average NOx emissions NBP-Participating States vs. the Control States 1997 - 2008. Middle: estimated time weights as of (2). Bottom: Point estimates and confidence intervals of DID and TWDID estimation.

untreated units, the TWDID estimator offers both a lower bias and a lower variance than the unweighted DID estimator. Moreover, I show how to adjust the standard errors to cover the additional weight estimation uncertainty. I revisit a study on the effect of a cap-and-trade program on NOx emissions. Using TWDID reduces the standard errors of the estimated average treatment effect by 10% compared to a conventional DID approach.



## References

- ABADIE, A., A. DIAMOND, AND J. HAINMUELLER (2010): “Synthetic control methods for comparative case studies: Estimating the effect of California’s tobacco control program,” *Journal of the American Statistical Association*, 105, 493–505.
- ABADIE, A. AND J. L’HOUR (2020): “A penalized synthetic control estimator for disaggregated data,” *Working Paper*.
- ARELLANO, M. (1987): “Computing robust standard errors for within-groups estimators,” *Oxford Bulletin of Economics and Statistics*, 49, 431–434.
- ARKHANGELSKY, D., S. ATHEY, D. A. HIRSHBERG, G. W. IMBENS, AND S. WAGER (2021): “Synthetic difference-in-differences,” *American Economic Review*, 111, 4088–4118.
- BAI, J. (2009): “Panel Data Models With Interactive Fixed Effects,” *Econometrica*, 77, 1229–1279.
- BERTRAND, M., E. DUFLO, AND S. MULLAINATHAN (2004): “How much should we trust differences-in-differences estimates?” *The Quarterly Journal of Economics*, 119, 249–275.
- CHAN, M. K. AND S. S. KWOK (2021): “The PCDID Approach: Difference-in-Differences When Trends Are Potentially Unparallel and Stochastic,” *Journal of Business & Economic Statistics*, 0, 1–18.
- DESCHENES, O., M. GREENSTONE, AND J. S. SHAPIRO (2017): “Defensive investments and the demand for air quality: Evidence from the NOx budget program,” *American Economic Review*, 107, 2958–89.
- FERMAN, B. (2019): “On the Properties of the Synthetic Control Estimator with Many Periods and Many Controls,” *arXiv preprint arXiv:1906.06665*.
- GOBILLON, L. AND T. MAGNAC (2016): “Regional policy evaluation: Interactive fixed effects and synthetic controls,” *Review of Economics and Statistics*, 98, 535–551.
- KETZ, P. (2018): “Subvector inference when the true parameter vector may be near or at the boundary,” *Journal of Econometrics*, 207, 285–306.
- PESARAN, M. H. (2006): “Estimation and Inference in Large Heterogeneous Panels with a Multifactor Error Structure,” *Econometrica*, 74, 967–1012.

## A Proofs of Theorems

### A.1 Proof of Theorem 1

Again some notation. Let  $\mathbf{b}$  be a  $T$ -vector with elements  $b_t = \frac{1}{T_1}I(t > T_0) - \frac{1}{T_0}I(t \leq T_0)$ .  $\mathbf{\Upsilon}$  is a  $T \times T$  diagonal matrix with the scaled time weights on the diagonal  $(\mathbf{\Upsilon})_{t,t} = T_0 v_t I(t \leq T_0) + I(t > T_0)$ . So  $\mathbf{b}'\mathbf{\Upsilon}\mathbf{x} = \frac{1}{T_1} \sum_{t>T_0} - \sum_{t \leq T_0} v_t x_t$  for any  $T$  vector  $\mathbf{x}$ . Let  $\bar{\nu}_\lambda = \bar{\lambda}^{(1)} - \bar{\lambda}^{(0)} - \boldsymbol{\xi}_\lambda$  and

$$Z_N(\mathbf{v}) = \mathbf{b}'\mathbf{\Upsilon}(\bar{\boldsymbol{\varepsilon}}^{(1)} - \bar{\boldsymbol{\varepsilon}}^{(0)}) = \frac{1}{T_1} \sum_{t>T_0} (\bar{\varepsilon}_t^{(1)} - \bar{\varepsilon}_t^{(0)}) - (\bar{\boldsymbol{\varepsilon}}_{[0]}^{(1)} - \bar{\boldsymbol{\varepsilon}}_{[0]}^{(0)})'\mathbf{v}$$

First write

$$\hat{\tau}(\mathbf{v}) - \tau = \boldsymbol{\xi}'_\lambda \boldsymbol{\xi}_f(\mathbf{v}) + \bar{\nu}'_\lambda \boldsymbol{\xi}_f(\mathbf{v}) + Z_N(\mathbf{v})$$

from which we arrive at the following asymptotic results.

LEMMA 1. For any fixed  $\mathbf{v} \in \mathbb{V}$

1.  $\sqrt{N}Z_N(\mathbf{v}) \xrightarrow{d} \mathbf{N}[0, V_\varepsilon(\mathbf{v})]$
2.  $\sqrt{N}(\bar{\lambda}^{(1)} - \bar{\lambda}^{(0)} - \boldsymbol{\xi}_\lambda) \xrightarrow{d} \mathbf{N}[0, \boldsymbol{\Sigma}_\lambda]$

with  $V_\varepsilon(\mathbf{v}) = \mathbf{b}'\mathbf{\Upsilon}\boldsymbol{\Sigma}_\varepsilon\mathbf{\Upsilon}\mathbf{b}$ ,  $\boldsymbol{\Sigma}_\varepsilon = \frac{\boldsymbol{\Sigma}_\varepsilon^{(0)}}{\kappa} + \frac{\boldsymbol{\Sigma}_\varepsilon^{(1)}}{1-\kappa}$  and  $\boldsymbol{\Sigma}_\lambda = \frac{\boldsymbol{\Sigma}_\lambda^{(0)}}{\kappa} + \frac{\boldsymbol{\Sigma}_\lambda^{(1)}}{1-\kappa}$

The first two assertions of Theorem 2 follow with

$$V_{\hat{\tau}}(\mathbf{v}) = \boldsymbol{\xi}_f(\mathbf{v})'\boldsymbol{\Sigma}_\lambda\boldsymbol{\xi}_f(\mathbf{v}) + V_\varepsilon(\mathbf{v})$$

### A.2 Proof of Theorem 2

The estimated time weights  $\hat{\mathbf{v}}$  as defined in (2) can be rewritten as

$$\hat{\mathbf{v}} = \arg \min_{\mathbf{v} \in \mathbb{V}} Q_N(\mathbf{v}), \quad Q_N(\mathbf{v}) = (\dot{\mathbf{y}}_{0,T} - \dot{\mathbf{Y}}_{0,p}\mathbf{v})'(\dot{\mathbf{y}}_{0,T} - \dot{\mathbf{Y}}_{0,p}\mathbf{v})$$

with  $\mathbf{y}_{0,t} = (y_{1t}, \dots, y_{N_0 t})$  the vector of the untreated units' outcomes in the period  $t$ ,  $\mathbf{Y}_{0,p} = [\mathbf{y}_{0,1}, \dots, \mathbf{y}_{0,T_0}]$  the  $N_0 \times T_0$  matrix of the untreated units' pre-treatment outcomes and  $\mathbb{V}$  the set of  $T_0$  vectors with non-zero elements that sum to one. We notice that  $\hat{\mathbf{v}}$  always exists, and is unique if  $\dot{\mathbf{Y}}_{0,p} = [\dot{\mathbf{y}}_{p,1}, \dots, \dot{\mathbf{y}}_{p,N_0}]'$  has full column rank.

Consistency follows from two key arguments. First, the objective function uniformly convergence to its population equivalent  $Q_\infty(\mathbf{v}) = V_f^{(0)}(\mathbf{v}) + V_\varepsilon(\mathbf{v})$

$$\sup_{\mathbf{v} \in \mathbb{V}} |Q_N(\mathbf{v}) - Q_\infty(\mathbf{v})| \xrightarrow{p} 0$$

Second,  $\mathbf{v}^*$  as defined in (5) is the uniquely identifiable minimizer of  $Q_\infty(\mathbf{v})$ .

To obtain asymptotic normality, suppose first that  $\mathbf{v}^*$  lies in the interior of  $\mathbb{V}$ , i.e. the non-negativity constraints are not binding in the limit. Without the non-negativity constraint we can rewrite the estimation as an unconstrained minimization and obtain an explicit solution. To do so, we transform the condition  $\sum_t v_t = 1$  into  $v_1 = 1 - \boldsymbol{\iota}'_{T_0-1} \mathbf{v}_{-1}$  with  $\mathbf{v}_{-1} = (v_2, \dots, v_{T_0})$ , so  $\mathbf{v} = \mathbf{e}_1 + \mathbf{R}\mathbf{v}_{-1}$  with  $\mathbf{R} = \begin{pmatrix} -\boldsymbol{\iota}'_{T_0-1} \\ \mathbf{I}_{T_0-1} \end{pmatrix}$ . The minimization problem then becomes

$$\min_{\mathbf{v}_{-1} \in \mathbb{R}^{T_0-1}} (\tilde{\mathbf{y}}_{(1)} - \tilde{\mathbf{Y}}_p \mathbf{v}_{-1})' (\tilde{\mathbf{y}}_{(1)} - \tilde{\mathbf{Y}}_p \mathbf{v}_{-1})$$

with  $\tilde{\mathbf{y}}_{(1)} = \bar{\mathbf{y}}_{(1)} - \mathbf{y}_1$  and  $\tilde{\mathbf{Y}}_p = \mathbf{Y}_p \mathbf{R}$ . The solution is

$$\hat{\mathbf{v}}_{-1} = (\tilde{\mathbf{Y}}_p' \tilde{\mathbf{Y}}_p)^{-1} \tilde{\mathbf{Y}}_p' \tilde{\mathbf{y}}_{(1)}$$

and the corresponding  $T_0$  vector of time weights will be  $\hat{\mathbf{v}} = (1 - \sum_{t=2}^{T_0} \hat{v}_t, \hat{\mathbf{v}}'_{-1})'$ . That is a least-squares regression with the first regressor subtracted. By the same argument we can write  $\mathbf{v}^* = \mathbf{e}_1 + \mathbf{R}\mathbf{v}_{-1}^*$ . The idea is to first show asymptotic normality for the last  $T_0 - 1$  weights,

$$\sqrt{N_0}(\hat{\mathbf{v}}_{-1} - \mathbf{v}_{-1}^*) \xrightarrow{d} \mathcal{N}[0, \boldsymbol{\Sigma}_{\hat{\mathbf{v}}_{-1}}]$$

which is easy to obtain because  $\hat{\mathbf{v}}_{-1}$  is an unrestricted least squares estimator. Because of the sum-to-1 condition it will immediately follow for the full vector

$$\sqrt{N_0}(\hat{\mathbf{v}} - \mathbf{v}^*) = \sqrt{N_0} \mathbf{R}(\hat{\mathbf{v}}_{-1} - \mathbf{v}_{-1}^*) \xrightarrow{d} \mathcal{N}[0, \boldsymbol{\Sigma}_{\hat{\mathbf{v}}}]$$

with  $\boldsymbol{\Sigma}_{\hat{\mathbf{v}}} = \mathbf{R} \boldsymbol{\Sigma}_{\hat{\mathbf{v}}_{-1}} \mathbf{R}'$ . Note that both  $\boldsymbol{\Sigma}_{\hat{\mathbf{v}}_{-1}}$  and  $\boldsymbol{\Sigma}_{\hat{\mathbf{v}}}$  have rank  $T_0 - 1$ . In absence of heteroskedasticity we suggest using the simple least squares variance estimator

$$\hat{\boldsymbol{\Sigma}}_{\hat{\mathbf{v}}_{-1}} = Q_N(\hat{\mathbf{v}}) (\mathbf{R}' \mathbf{Y}'_p \mathbf{Y}_p \mathbf{R})^{-1}$$

Suppose now that  $\mathbf{v}^*$  lies on the boundary of  $\mathbb{V}$ . That is, at least one element is exactly zero and  $\mathbf{v}^*$  is sparse. In general, asymptotic normality of extremum estimators can break down near the boundary of the parameter space (Ketz, 2018). In this special case, however, the part of  $\hat{\mathbf{v}}$  belonging to the non-zero elements of  $\mathbf{v}^*$  will still be asymptotically normal.

Formally, suppose without loss of generality that the first  $k$  elements ( $0 < k < T_0 - 1$ ) of  $\mathbf{v}^*$  are zero and write  $\mathbf{v} = (\mathbf{v}'_k, \mathbf{v}'_{-k})'$ .

LEMMA 2. *I conjecture:  $\hat{\mathbf{v}}_k = O_p(N^{-q})$  with  $q > 0.5$*

The idea is that  $\hat{\mathbf{v}}_k$  converges to zero so quickly that the asymptotic distribution is the same as if we had just set  $\hat{\mathbf{v}}_k = 0$  in the first place. But this means that the first  $k$  periods are irrelevant and we have reduced the exercise to  $T_0 - k$  pre-treatment periods, for which the pseudo-true weights  $\mathbf{v}_{-k}^*$  are not sparse. Then we can apply the results from the previous paragraph to get asymptotic normality for the non-zero elements.

### A.3 Proof of Theorem 3

Let  $\mathbf{G}_N = (\bar{\boldsymbol{\nu}}'_\lambda \boldsymbol{\xi}_f(\mathbf{v}^*) + Z_N(\mathbf{v}^*), (\hat{\mathbf{v}} - \mathbf{v}^*)' )'$ .

LEMMA 3. 1.  $(\mathbf{F}'_{(0)} \bar{\boldsymbol{\nu}}_\lambda)'(\hat{\mathbf{v}} - \mathbf{v}^*) = O_p(\frac{1}{N})$

2.  $(\bar{\boldsymbol{\varepsilon}}^{(1)}_{(0)} - \bar{\boldsymbol{\varepsilon}}^{(0)}_{(0)})'(\hat{\mathbf{v}} - \mathbf{v}^*) = O_p(\frac{1}{N})$

3.  $\sqrt{N}G_N \xrightarrow{d} N[\mathbf{0}, \boldsymbol{\Sigma}_G]$

with

$$\boldsymbol{\Sigma}_G = \begin{pmatrix} V_{\hat{\tau}} & \frac{1}{\sqrt{\kappa}} \mathbf{C}' \\ \frac{1}{\sqrt{\kappa}} \mathbf{C} & \frac{1}{\kappa} \boldsymbol{\Sigma}_{\hat{\mathbf{v}}} \end{pmatrix}$$

Adding and subtracting  $\hat{\tau}(\mathbf{v}^*)$  we get

$$\hat{\tau}(\hat{\mathbf{v}}) - \tau - b(\mathbf{v}^*) = (1, -(\mathbf{F}_{(0)} \boldsymbol{\xi}_\lambda)') \mathbf{G}_N + O_p(\frac{1}{N})$$

and thus assertions 1 and 2 of Theorem 3 follow with

$$V_{\hat{\tau}(\hat{\mathbf{v}})} = V_{\hat{\tau}} + \frac{1}{\kappa} (\mathbf{F}_{(0)} \boldsymbol{\xi}_\lambda)' \boldsymbol{\Sigma}_{\hat{\mathbf{v}}} \mathbf{F}_{(0)} \boldsymbol{\xi}_\lambda - \frac{2}{\sqrt{\kappa}} (\mathbf{F}_{(0)} \boldsymbol{\xi}_\lambda)' \mathbf{C}$$

Consistency of  $\widehat{V}_{\hat{\tau}(\hat{\mathbf{v}})}$  follows from the fact that  $\widehat{V}_{\hat{\tau}}(\hat{\mathbf{v}}) \xrightarrow{p} V_{\hat{\tau}}(\mathbf{v}^*)$  and  $\dot{\mathbf{A}}_{(0)} \xrightarrow{p} \mathbf{F}_{(0)} \boldsymbol{\xi}_\lambda$ . So far  $\widehat{V}_{\hat{\tau}(\hat{\mathbf{v}})}$  ignores the covariance term  $\mathbf{C}$ . In simulations  $\mathbf{C}$  is negligible as long as  $\hat{\mathbf{v}}$  is estimated using demeaned observations.

### A.4 Proof of Theorem 4

Let  $\dot{\mathbf{y}}_i = \mathbf{y}_i - \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i$  denote the  $T$ -vector of cross-sectionally demeaned outcomes,  $\ddot{\mathbf{y}}_i = \mathbf{M}_T \dot{\mathbf{y}}_i$  is the vector of double demeaned outcomes with  $\mathbf{M}_T = \mathbf{I}_T - \frac{\iota_T \iota_T'}{T}$ .

I now derive probability limit of  $\widehat{V}(\mathbf{v})$  for some fixed  $\mathbf{v} \in \mathbb{V}$ . First note that the denominator reduces to  $\frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \ddot{\mathbf{D}}_i = T k_N k_T (1 - k_N)(1 - k_T) =: Q_N$ . Hence we only need to consider the numerator.

$$\begin{aligned} \hat{\mathbf{u}}_{i,\mathbf{v}} &= \ddot{y}_{i,\mathbf{v}} - \ddot{\mathbf{D}}_i \hat{\tau}(\mathbf{v}) \\ \ddot{y}_{i,\mathbf{v}} &= \mathbf{M}_T \boldsymbol{\Upsilon} \dot{\mathbf{y}}_i = \tau \mathbf{M}_T \boldsymbol{\Upsilon} \dot{\mathbf{D}}_i + \ddot{\mathbf{u}}_{i,\mathbf{v}} \\ \ddot{\mathbf{u}}_{i,\mathbf{v}} &= \mathbf{M}_T (\mathbf{F}_v \dot{\boldsymbol{\lambda}}_i + \dot{\boldsymbol{\varepsilon}}_{i,\mathbf{v}}) \end{aligned}$$

with  $\mathbf{F}_v = \boldsymbol{\Upsilon} \mathbf{F}$  and  $\dot{\boldsymbol{\varepsilon}}_{i,\mathbf{v}} = \boldsymbol{\Upsilon} \dot{\boldsymbol{\varepsilon}}_i$ . Note that  $\boldsymbol{\Upsilon} \dot{\mathbf{D}}_i = \dot{\mathbf{D}}_i$  because  $(\boldsymbol{\Upsilon})_{t,t} \neq 1$

only if  $D_{it} = 0$ . Substitute  $\hat{\mathbf{u}}_{i,v} = \ddot{\mathbf{u}}_{i,v} - (\hat{\tau} - \tau)\ddot{\mathbf{D}}_i$  to get

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \hat{\mathbf{u}}_{i,v} \hat{\mathbf{u}}_{i,v}' \ddot{\mathbf{D}}_i &= \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \ddot{\mathbf{u}}_{i,v} \ddot{\mathbf{u}}_{i,v}' \ddot{\mathbf{D}}_i \\ &\quad - 2(\hat{\tau}(\mathbf{v}) - \tau) \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \ddot{\mathbf{u}}_{i,v} \ddot{\mathbf{D}}_i \ddot{\mathbf{D}}_i \\ &\quad + (\hat{\tau}(\mathbf{v}) - \tau)^2 \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \ddot{\mathbf{D}}_i \ddot{\mathbf{D}}_i \ddot{\mathbf{D}}_i \end{aligned}$$

The following lemma ensures that the latter two terms converge to zero given consistency of  $\hat{\tau}(\mathbf{v})$ .

LEMMA 4. *We have these convergence results as  $N \rightarrow \infty$*

- $\frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \ddot{\mathbf{u}}_{i,v} \ddot{\mathbf{D}}_i \ddot{\mathbf{D}}_i \xrightarrow{p} c_1 < \infty$
- $\frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \ddot{\mathbf{D}}_i \ddot{\mathbf{D}}_i \ddot{\mathbf{D}}_i \xrightarrow{p} c_2 < \infty$
- $\frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \mathbf{F}_v \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\varepsilon}}_{i,v}' \ddot{\mathbf{D}}_i \xrightarrow{p} 0$

Use  $\ddot{\mathbf{D}}_i' \ddot{\mathbf{u}}_{i,v} = \ddot{\mathbf{D}}_i' \mathbf{M}_T' \mathbf{M}_T \ddot{\mathbf{u}}_{i,v} = \ddot{\mathbf{D}}_i' \ddot{\mathbf{u}}_{i,v}$ ,  $\ddot{\mathbf{u}}_{i,v} = \mathbf{F}_v \dot{\boldsymbol{\lambda}}_i + \dot{\boldsymbol{\varepsilon}}_{i,v}$  to get

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \ddot{\mathbf{u}}_{i,v} \ddot{\mathbf{u}}_{i,v}' \ddot{\mathbf{D}}_i &= \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \boldsymbol{\Upsilon} \mathbf{F} \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\lambda}}_i' \mathbf{F}' \boldsymbol{\Upsilon} \ddot{\mathbf{D}}_i + \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \dot{\boldsymbol{\varepsilon}}_{i,v} \dot{\boldsymbol{\varepsilon}}_{i,v}' \ddot{\mathbf{D}}_i \quad (8) \\ &\quad + \frac{2}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \mathbf{F}_v \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\varepsilon}}_{i,v}' \ddot{\mathbf{D}}_i \end{aligned}$$

where the last term converges to zero by independence of  $\mathbf{F}$  and  $\boldsymbol{\varepsilon}_i$ . Consider the first part of (8).

Then we can write  $\ddot{\mathbf{D}}_i = Q_N \mathbf{b}(\frac{1}{1-k_N} I(i > N_0) - \frac{1}{k_N} I(i \leq N_0))$ . Also,  $\mathbf{b}' \boldsymbol{\Upsilon} \mathbf{F} = \bar{\mathbf{f}}_{(1)} - \mathbf{F}'_p \mathbf{v}$ . Let  $\Sigma_\lambda^{(0)} := \frac{1}{N_0} \sum_{i=1}^{N_0} \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\lambda}}_i'$  and  $\Sigma_\lambda^{(1)} := \frac{1}{N_1} \sum_{i=N_0+1}^N \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\lambda}}_i'$ . Then, all together we get

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \boldsymbol{\Upsilon} \mathbf{F} \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\lambda}}_i' \mathbf{F}' \boldsymbol{\Upsilon} \ddot{\mathbf{D}}_i &= \frac{Q_N^2}{N k_N^2} \sum_{i=1}^{N_0} \mathbf{b}' \boldsymbol{\Upsilon} \mathbf{F} \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\lambda}}_i' \mathbf{F}' \boldsymbol{\Upsilon} \mathbf{b} \\ &\quad + \frac{Q_N^2}{N(1-k_N)^2} \sum_{i=N_0}^N \mathbf{b}' \boldsymbol{\Upsilon} \mathbf{F} \dot{\boldsymbol{\lambda}}_i \dot{\boldsymbol{\lambda}}_i' \mathbf{F}' \boldsymbol{\Upsilon} \mathbf{b} \\ &= Q_N^2 \left[ (\bar{\mathbf{f}}_{(1)} - \mathbf{F}'_p \mathbf{v})' \left( \frac{\Sigma_\lambda^{(0)}}{k_N} + \frac{\Sigma_\lambda^{(1)}}{1-k_N} \right) (\bar{\mathbf{f}}_{(1)} - \mathbf{F}'_p \mathbf{v}) \right] \end{aligned}$$

Consider now the second part of (8). First, observe that

$$\frac{\frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \dot{\boldsymbol{\varepsilon}}_{i,v}}{Q_N} = \mathbf{b}' \boldsymbol{\Upsilon} \boldsymbol{\Delta}_\varepsilon = \frac{1}{T_1} \sum_{t>T_0} \Delta_{\varepsilon,t} - \sum_{t \leq T_0} v_t \Delta_{\varepsilon,t}$$

with  $\mathbf{\Delta}_\varepsilon = \bar{\boldsymbol{\varepsilon}}^{(1)} - \bar{\boldsymbol{\varepsilon}}^{(0)}$ . We should therefore have that

$$\frac{\frac{1}{N} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \dot{\boldsymbol{\varepsilon}}_{i,\mathbf{v}} \dot{\boldsymbol{\varepsilon}}_{i,\mathbf{v}}' \ddot{\mathbf{D}}_i}{Q_N^2} \xrightarrow{p} \frac{\text{var}[\frac{1}{\sqrt{N}} \sum_{i=1}^N \ddot{\mathbf{D}}_i' \dot{\boldsymbol{\varepsilon}}_{i,\mathbf{v}}]}{Q_N^2} = N \text{ var}[\mathbf{b}' \boldsymbol{\Upsilon} \mathbf{\Delta}_\varepsilon].$$

Moreover, by independence of the units we have  $N \text{ var}[\mathbf{\Delta}_\varepsilon] = N \text{ var}[\bar{\boldsymbol{\varepsilon}}^{(1)}] + N \text{ var}[\bar{\boldsymbol{\varepsilon}}^{(0)}] = \frac{\boldsymbol{\Omega}}{k_N(1-k_N)}$  and thus

$$V_\varepsilon(\mathbf{v}) = \frac{\mathbf{b}' \boldsymbol{\Upsilon} \boldsymbol{\Omega} \boldsymbol{\Upsilon} \mathbf{b}}{k_N(1-k_N)} = \frac{\bar{\sigma}_{(0),\mathbf{v}} + \bar{\sigma}_{(1)} - 2\bar{\sigma}_{(0,1),\mathbf{v}}}{k_N(1-k_N)}$$

with  $\bar{\sigma}_{(0),\mathbf{v}} = \mathbf{v}' \boldsymbol{\Omega}_{(0)} \mathbf{v}$ ,  $\bar{\sigma}_{(1)} = \frac{1}{T_1^2} \boldsymbol{\nu}'_{T_1} \boldsymbol{\Omega}_{(1)} \boldsymbol{\nu}_{T_1}$  and  $\bar{\sigma}_{(0,1),\mathbf{v}} = \frac{1}{T_1} \mathbf{v}' \boldsymbol{\Omega}_{(0,1)} \boldsymbol{\nu}_{T_1}$ .