

NBER WORKING PAPER SERIES

TWO-WAY FIXED EFFECTS AND DIFFERENCES-IN-DIFFERENCES ESTIMATORS
WITH SEVERAL TREATMENTS

Clément de Chaisemartin
Xavier D'Haultfoeulle

Working Paper 30564
<http://www.nber.org/papers/w30564>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
October 2022

Several of this paper's ideas arose during conversations with Enrico Cantoni, Angelica Meinhofer, Vincent Pons, Jimena Rico-Straffon, Marc Sangnier, Oliver Vanden Eynde, and Liam Wren-Lewis who shared with us their interrogations, and sometimes their referees' interrogations, on two-way fixed effects regressions with several treatments. We are grateful to them for those stimulating conversations. We are grateful to Yubo Wei for his excellent work as a research assistant. We are also grateful to the editor, associate editor, and two anonymous referees for their very helpful comments. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2022 by Clément de Chaisemartin and Xavier D'Haultfoeulle. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Two-way Fixed Effects and Differences-in-Differences Estimators with Several Treatments
Clément de Chaisemartin and Xavier D'Haultfoeuille
NBER Working Paper No. 30564
October 2022
JEL No. C21,C23

ABSTRACT

We study two-way-fixed-effects regressions (TWFE) with several treatment variables. Under a parallel trends assumption, we show that the coefficient on each treatment identifies a weighted sum of that treatment's effect, with possibly negative weights, plus a weighted sum of the effects of the other treatments. Thus, those estimators are not robust to heterogeneous effects and may be contaminated by other treatments' effects. When a treatment is omitted from the regression, we obtain a new omitted variable bias formula, where bias can arise even if the treatments are not correlated with each other, but can be smaller than in the TWFE regression with all treatments. We propose an alternative difference-in-differences estimator, robust to heterogeneous effects and immune to the contamination problem. In the application we consider, the TWFE regression identifies a highly non-convex combination of effects, with large contamination weights, and one of its coefficients significantly differs from our heterogeneity-robust estimator.

Clément de Chaisemartin
Department of Economics
University of California at Santa Barbara
Santa Barbara, CA 93106
and NBER
clementdechaisemartin@ucsb.edu

Xavier D'Haultfoeuille
CREST
5 avenue Henry Le Chatelier
91764 Palaiseau cedex
FRANCE
xavier.dhaultfoeuille@ensae.fr

Two-way Fixed Effects and Differences-in-Differences Estimators with Several Treatments*

Clément de Chaisemartin

Xavier D’Haultfoeuille[†]

Abstract

We study two-way-fixed-effects regressions (TWFE) with several treatment variables. Under a parallel trends assumption, we show that the coefficient on each treatment identifies a weighted sum of that treatment’s effect, with possibly negative weights, plus a weighted sum of the effects of the other treatments. Thus, those estimators are not robust to heterogeneous effects and may be contaminated by other treatments’ effects. When a treatment is omitted from the regression, we obtain a new omitted variable bias formula, where bias can arise even if the treatments are not correlated with each other, but can be smaller than in the TWFE regression with all treatments. We propose an alternative difference-in-differences estimator, robust to heterogeneous effects and immune to the contamination problem. In the application we consider, the TWFE regression identifies a highly non-convex combination of effects, with large contamination weights, and one of its coefficients significantly differs from our heterogeneity-robust estimator.

(*JEL C21, C23*)

1 Introduction

To estimate treatment effects, researchers often use panels of groups (e.g. counties, regions), and estimate two-way fixed effect (TWFE) regressions, namely regressions of the outcome variable on group and time fixed effects and the treatment. de Chaisemartin and D’Haultfoeuille (2020)

*Several of this paper’s ideas arose during conversations with Enrico Cantoni, Angelica Meinhofer, Vincent Pons, Jimena Rico-Straffon, Marc Sangnier, Oliver Vanden Eynde, and Liam Wren-Lewis who shared with us their interrogations, and sometimes their referees’ interrogations, on two-way fixed effects regressions with several treatments. We are grateful to them for those stimulating conversations. We are grateful to Yubo Wei for his excellent work as a research assistant. We are also grateful to the editor, associate editor, and two anonymous referees for their very helpful comments.

[†]de Chaisemartin: Sciences Po (email: clement.dechaisemartin@sciencespo.fr); D’Haultfoeuille: CREST-ENSAE (email: xavier.dhaultfoeuille@ensae.fr). Xavier D’Haultfoeuille thanks the hospitality of PSE where this research was conducted.

have found that almost 20% of empirical papers published by the American Economic Review (AER) from 2010 to 2012 estimate such regressions.

Under a parallel trends assumption, TWFE regressions with one treatment identify a weighted sum of the treatment effects of treated (g, t) cells, with weights that may be negative and sum to one (see de Chaisemartin and D’Haultfœuille, 2020; Borusyak and Jaravel, 2017). Because of the negative weights, the treatment coefficient in such regressions is not robust to heterogeneous treatment effects across groups and time periods: it may be, say, negative, even if the treatment effect is strictly positive in every (g, t) cell.

However, in 18% of the TWFE papers published in the AER from 2010 to 2012, the TWFE regression has several treatment variables. By including several treatments, researchers hope to estimate the effect of each treatment holding the other treatments constant. For instance, when studying the effect of marijuana laws, as in Meinhofer et al. (2021), one may want to separate the effect of medical and recreational laws. To do so, one may estimate a regression of the outcome of interest in state g and year t on state fixed effects, year fixed effects, an indicator for whether state g has a medical law in year t , and an indicator for whether state g has a recreational law in year t .

In this paper, we investigate what TWFE regressions with several treatments identify. We show that under a parallel trends assumption, the coefficient on each treatment identifies the sum of two terms. The first term is a weighted sum of the effect of that treatment in each group and period, with weights that may be negative and sum to one. A similar weighted sum appears in decompositions of TWFE regressions with only one treatment. The second term is a sum of the effects of the other treatments, with weights summing to zero. Accordingly, with several treatments, coefficients in TWFE regressions may be contaminated by the effect of other treatments, an issue that was not present with one treatment. As the weights sum to zero, this second term disappears if the effect of the other treatments is homogeneous, but it is often implausible that those effects are homogeneous. The weights attached to any TWFE regression with several treatments can be computed by the `twowayfeweights` Stata and R packages. Estimating those weights may be useful, to assess if a TWFE coefficient is robust to heterogeneous treatment effects, and if it is contaminated by the effect of the other treatments in the regression.

We consider simple examples with two treatments, to show that TWFE regressions may not be robust to heterogeneous effects because they may leverage two types of “forbidden comparisons”, borrowing the terminology coined by Borusyak and Jaravel (2017). In a first example, the coefficient on the first treatment leverages a difference-in-differences (DID) comparing the outcome evolution of a group going from untreated to receiving both treatments to the outcome evolution of a “control” group going from untreated to receiving the second treatment. If the effect of the second treatment is the same in the two groups, those two effects cancel each other

out in this DID. But if the effects of the second treatment differ in the two groups, they do not cancel each other out, and they contaminate the coefficient on the first treatment. In a second example, the coefficient on the first treatment leverages a DID comparing the outcome evolution of a group going from untreated to receiving the first treatment to the outcome evolution of a “control” group that receives the second treatment at both periods. If the control group’s effect of the second treatment is the same in the pre and in the post period, those two effects cancel each other out in this DID. But if the control group’s effect of the second treatment changes over time, those two effects do not cancel out, and they contaminate the coefficient on the first treatment.

We then consider a TWFE regression that would omit the other treatments, and derive an analogue of the standard omitted variable bias (OVB) formula for that regression, allowing for heterogeneous effects. Allowing for heterogeneous effects leads to two perhaps surprising departures from the standard OVB logic many researchers are used to. First, omitting from a TWFE regression a treatment uncorrelated with the main treatment of interest may still lead to an OVB. Second, controlling for more treatments may lead to a more biased estimator than not controlling for them. Specifically, we show that in the presence of two treatments, a TWFE regression with only the first treatment also estimates a weighted sum of the effect of that treatment in each group and period, with weights that may be negative and sum to one, plus a weighted sum of the effects of the other treatments, but with weights that do not sum to zero. That second term, which corresponds to the OVB term in the standard formula, may differ from zero even if the two treatments are uncorrelated conditional on the group and time fixed effects. Then, we use our decompositions of the TWFE regressions with one and several treatments to derive the maximal bias of both regressions for the average effect of the first treatment on the treated, under the assumption that the effect of every treatment is bounded in absolute value by a (potentially large) constant in every group and period. The ratio between the maximal biases of both regressions is independent of that constant and can be estimated, thus allowing researchers to compare the maximal bias of the two regressions. The ratio of the regressions’ maximal biases can either be smaller or larger than one in practice.

Finally, we propose an alternative DID estimator that relies on common trends assumptions, like TWFE regressions, but that is robust to heterogeneous effects and does not suffer from the contamination problem, unlike TWFE regressions. Our estimator generalizes the DID_M estimator in de Chaisemartin and D’Haultfœuille (2020) to instances with several treatments. To isolate the effect of the first treatment, our estimator compares the $t - 1$ -to- t outcome evolution, of switching groups whose first treatment switches from $t - 1$ to t while their other treatments do not change, and of control groups i) whose treatments all remain the same, and ii) that had the same treatments as the switching groups in period $t - 1$. i) ensures that our new estimator is robust to heterogeneous effects across groups of all treatments. ii) ensures that it is robust to

heterogeneous effects over time of all treatments.

Our estimator’s robustness may come at a high price in terms of external validity and statistical precision. For instance, in our application in Section 5, we can only match a small number of switchers to valid control groups meeting i) and ii). Then, there may be internal-external validity and bias-variance trade-offs between our new estimator and less robust estimators, such as the DID_M estimator in de Chaisemartin and D’Haultfœuille (2020) or TWFE regressions with several treatments. To account for the fact our new estimator may sometimes be estimated on a small sample of groups, we propose, in addition to a standard confidence interval that is asymptotically valid under weak conditions, another confidence interval that has both exact coverage under a normality assumption and is asymptotically valid without such a normality requirement.

As an illustration, we use our results to revisit Hotz and Xiao (2011), who run TWFE regressions of measures of daycare quality in state g and year t on two daycare regulations in state g and year t : the minimum number of years of schooling required to be a daycare director and the minimum staff-to-child ratio. Focusing on the years-of-schooling treatment, we find that the TWFE regression with several treatments estimates weighted sums of effects with very large negative weights attached to them, both on the treatment’s own effects, but also on the effects of the other treatments in the regression. The TWFE regression with only the years-of-schooling treatment has much smaller weights attached to it. As a result, the maximal bias of the TWFE regression with several treatments is almost five times larger than that of the regression including only the years-of-schooling treatment. Thus, the “short” regression seems preferable, at least per our maximal-bias metric. We finally show that our heterogeneity-robust estimator is much closer to zero than, and significantly different from, the coefficient of the TWFE regression with several treatments.

Our paper is closely related to the recent literature showing that TWFE regressions with one treatment variable may not be robust to heterogeneous effects (see de Chaisemartin and D’Haultfœuille, 2020; Goodman-Bacon, 2021; Borusyak and Jaravel, 2017). Our paper is also closely related to several papers that have considered the causal interpretation of OLS regression coefficients with several treatments (see Sun and Abraham, 2021; Hull, 2018; Goldsmith-Pinkham, Hull and Kolesár, 2021). We discuss those papers in more details later in the paper (see Section 3.4), but for now we just note that when the treatments are indicators for whether group g has started receiving a binary and staggered treatment ℓ periods ago, our decomposition of the TWFE regression reduces to one of the decompositions in Sun and Abraham (2021). Accordingly, our decomposition extends their result to situations where the treatments in the regression are different policies that could non-binary, non-staggered, and non-mutually-exclusive, rather than indicators for having received a single binary and staggered policy ℓ periods ago.

The remainder of the paper is organized as follows. Section 2 presents the set up. Section 3

presents our decomposition results for TWFE regressions with several treatments. Section 4 presents our alternative estimator. Section 5 presents our empirical application.

2 Set up

We consider a panel of G groups observed at T periods, respectively indexed by g and t . Typically, groups are geographical entities gathering many observations, but a group could also just be a single individual or firm. For every $(g, t) \in \{1, \dots, G\} \times \{1, \dots, T\}$, let $N_{g,t}$ denote the population of cell (g, t) , and let $N = \sum_{g,t} N_{g,t}$ be the total population across all cells.

We are interested in the effect of K treatments. In this paper, we follow, e.g., Holland (1986); Holland and Rubin (1987), and define as a treatment a variable that has a causal effect on the outcome, in the sense that different values of that variable lead to different counterfactual outcomes. For every $(k, g, t) \in \{1, \dots, K\} \times \{1, \dots, G\} \times \{1, \dots, T\}$, let $D_{g,t}^k$ denote the value of treatment k for group g at period t , and let $D_{g,t} = (D_{g,t}^k)_{k \in \{1, \dots, K\}}$ denote a vector stacking together the K treatments of group g at period t . For every k , let \mathcal{D}_k denote the values $D_{g,t}^k$ can take. For now, we assume that the treatments are binary: $\mathcal{D}_k = \{0, 1\}$ for all k . This is just to simplify the exposition: our results can be extended to non-binary treatments, as explained below. For any $\mathbf{d} \in \{0, 1\}^K$, let $Y_{g,t}(\mathbf{d})$ denote the potential outcome of group g at period t if $(D_{g,t}^1, \dots, D_{g,t}^K) = \mathbf{d}$. The observed outcome is $Y_{g,t} = Y_{g,t}(D_{g,t}^1, \dots, D_{g,t}^K)$.

Importantly, our notation does not necessarily rule out dynamic effects of past treatments on the outcome. The K treatments may for instance include lags of the same treatment variables. We discuss this issue in more details after Theorem 2 below, and in Section 4.4.

We consider the treatments and potential outcomes of each (g, t) cell as random variables. For instance, aggregate random shocks may affect the potential outcomes of group g at period t , and that cell's treatments may also be random. All expectations below are taken with respect to the distribution of those random variables. On the other hand, the populations of cells (g, t) $N_{g,t}$ are treated as non-random throughout the paper.

Throughout the paper, we maintain the following assumptions. Below, we let $\mathbf{0} = (0, \dots, 0)$ denote the vector of K zeros.

Assumption 1 (*Balanced panel of groups*) For all $(g, t) \in \{1, \dots, G\} \times \{1, \dots, T\}$, $N_{g,t} > 0$.

Assumption 2 (*Independent groups*) The vectors $((Y_{g,t}(\mathbf{d}))_{\mathbf{d} \in \{0,1\}^K}, (D_{g,t}^k)_{k \in \{1, \dots, K\}})_{t \in \{1, \dots, T\}}$ are mutually independent.

Assumption 3 (*Strong exogeneity and common trends*) For all $(g, t) \in \{1, \dots, G\} \times \{2, \dots, T\}$,

1. $E(Y_{g,t}(\mathbf{0}) - Y_{g,t-1}(\mathbf{0}) | D_{g,1}, \dots, D_{g,T}) = E(Y_{g,t}(\mathbf{0}) - Y_{g,t-1}(\mathbf{0}))$.

2. $E(Y_{g,t}(\mathbf{0}) - Y_{g,t-1}(\mathbf{0}))$ does not vary across g .

Assumption 1 requires that no group appears or disappears over time. Assumption 2 requires that potential outcomes and treatments of different groups be independent, but it allows these variables to be correlated over time within each group. This is a commonly-made assumption in DID analysis, where standard errors are usually clustered at the group level (see Bertrand, Duflo and Mullainathan, 2004). Point 1 of Assumption 3 is related to the strong exogeneity condition in panel data models. It requires that the shocks affecting group g 's untreated outcome be mean independent of group g 's treatments. For instance, this rules out cases where a group gets treated because it experiences negative shocks, the so-called Ashenfelter's dip (see Ashenfelter, 1978). Point 2 requires that in every group, the expectation of the untreated outcome follow the same evolution over time. It is a generalization of the standard common trends assumption in DID models (see, e.g., Abadie, 2005).

We now define the TWFE regression described in the introduction, as well as our estimand of interest β_{fe} , the expectation of the treatment coefficient in the regression.¹

Regression 1 (*TWFE regression with K treatments*)

Let $\beta_{fe} = E[\hat{\beta}_{fe}]$, where $\hat{\beta}_{fe}$ denotes the coefficient on $D_{g,t}^1$ in a sample OLS regression of $Y_{g,t}$ on group fixed effects, period fixed effects, and the vector $D_{g,t}$, weighted by $N_{g,t}$.²

$$Y_{g,t} = \hat{\alpha}_g + \hat{\gamma}_t + \hat{\beta}_{fe} D_{g,t} + u_{g,t}, \quad (1)$$

where $u_{g,t}$ denotes the regression residual.

On top of the K treatments, the regression may also include some covariates. The decompositions below can easily be extended to this case, following the same steps as those used by de Chaisemartin and D'Haultfoeuille (2020) to extend their decomposition of TWFE regressions with one treatment to TWFE regressions with one treatment and some covariates (see Theorem S4 therein). In Section 3.4.2 below, we elaborate on the difference between a treatment and a covariate.

Let \mathbf{D} be the vector $(D_{g,t})_{(g,t) \in \{1, \dots, G\} \times \{1, \dots, T\}}$ collecting all the treatments in all the (g, t) cells. let $\mathbf{D}_g = (D_{1,g}, \dots, D_{T,g})$ be the vector collecting all the treatments in group g . Let $N_1 = \sum_{g,t} N_{g,t} D_{g,t}^1$ denote the total population of cells receiving the first treatment. Let $D_{g,t}^{-1} = (D_{g,t}^2, \dots, D_{g,t}^K)$ denote

¹ Throughout the paper, we assume that the treatments $D_{g,t}^k$ in Regression 1 are not collinear with the other independent variables in those regressions, so $\hat{\beta}_{fe}$ is well-defined.

² The regression could also be estimated using more disaggregated outcome data. For instance, groups may be US counties, and one may estimate the regression using individual-level outcome measures. This disaggregated regression is equivalent to the aggregated regression in (1), provided $Y_{g,t}$ is defined as the average outcome of individuals in cell (g, t) , and the aggregated regression is weighted by the number of individuals in cell (g, t) . Accordingly, the results below also apply to disaggregated regressions.

a vector stacking together the treatments of cell (g, t) , excluding treatment 1. Let $\varepsilon_{g,t}$ denote the residual of cell (g, t) in the sample regression of $D_{g,t}^1$ on group and period fixed effects and $D_{g,t}^{-1}$:

$$D_{g,t}^1 = \hat{\alpha} + \hat{\gamma}_g + \hat{\nu}_t + (D_{g,t}^{-1})' \hat{\zeta} + \varepsilon_{g,t}. \quad (2)$$

One can show that if the regressors in Regression 1 are not collinear, the average value of $\varepsilon_{g,t}$ across all (g, t) cells with $D_{g,t}^1 = 1$ differs from 0: $\sum_{(g,t):D_{g,t}^1=1} (N_{g,t}/N_1) \varepsilon_{g,t} \neq 0$. Then we let $w_{g,t}$ denote $\varepsilon_{g,t}$ divided by that average:

$$w_{g,t} = \frac{\varepsilon_{g,t}}{\sum_{(g,t):D_{g,t}^1=1} (N_{g,t}/N_1) \varepsilon_{g,t}}.$$

3 TWFE regressions with several treatments and heterogeneous effects

3.1 Decomposition results

3.1.1 Two treatment variables

For expositional purposes, we begin by considering the case with two treatments. For any $(g, t) \in \{1, \dots, G\} \times \{1, \dots, T\}$, let

$$\Delta_{g,t}^2 = Y_{g,t}(0, 1) - Y_{g,t}(0, 0)$$

denote the effect, in cell (g, t) , of moving the second treatment from zero to 1 while keeping the first treatment at zero. Let also

$$\Delta_{g,t}^1 = Y_{g,t}(1, D_{g,t}^2) - Y_{g,t}(0, D_{g,t}^2)$$

denote the effect, in cell (g, t) , of moving the first treatment from zero to one while keeping the second treatment at its observed value. When one estimates a TWFE regression with two treatments, a natural target parameter for β_{fe} , the coefficient on the first treatment, is

$$\delta_{ATT} = E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} \Delta_{g,t}^1 \right],$$

the average effect of moving $D_{g,t}^1$ from 0 to 1 while keeping $D_{g,t}^2$ at its observed value, across all (g, t) s such that $D_{g,t}^1 = 1$. δ_{ATT} is the ATT of $D_{g,t}^1$ controlling for $D_{g,t}^2$. We now show that β_{fe} does not identify δ_{ATT} in general.

Theorem 1 *Suppose that Assumptions 1-3 hold and $K = 2$. Then,*

$$\beta_{fe} = E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t}^1 + \sum_{(g,t):D_{g,t}^2=1} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t}^2 \right]. \quad (3)$$

Moreover, $\sum_{(g,t):D_{g,t}^1=1}(N_{g,t}/N_1)w_{g,t} = 1$ and $\sum_{(g,t):D_{g,t}^2=1}(N_{g,t}/N_1)w_{g,t} = 0$.

Theorem 1 shows that the coefficient on $D_{g,t}^1$ identifies the sum of two terms. The first term is a weighted sum of the average effect of moving $D_{g,t}^1$ from 0 to 1 while keeping $D_{g,t}^2$ at its observed value, across all (g, t) such that $D_{g,t}^1 = 1$, and with weights summing to 1. The second term is a weighted sum of the effect of moving $D_{g,t}^2$ from 0 to 1 while keeping $D_{g,t}^1$ at 0, across all (g, t) such that $D_{g,t}^2 = 1$, and with weights summing to 0. If the effect of $D_{g,t}^2$ is constant ($\Delta_{g,t}^2 = \delta^2$ for all (g, t)), this second term is equal to zero, but it may differ from zero if the effect of $D_{g,t}^2$ is heterogeneous.

Theorem 1 implies that there are two reasons why β_{fe} may differ from δ_{ATT} . First, some of the weights $w_{g,t}$ may differ from one. When the weights $w_{g,t}$ differ from one, one may have that

$$E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t}^1 \right] \neq \delta_{ATT},$$

if the effect of $D_{g,t}^1$ is heterogeneous across (g, t) cells. Some of the weights $w_{g,t}$ could even be negative, in which case $E \left[\sum_{(g,t):D_{g,t}^1=1} (N_{g,t}/N_1) w_{g,t} \Delta_{g,t}^1 \right]$ does not satisfy the no-sign reversal property: this quantity could for instance be negative, even if $\Delta_{g,t}^1 \geq 0$ for all (g, t) . With two treatments, negative weights can occur even in very simple designs, where there would not be any negative weights in the absence of the second treatment. For instance, consider a standard DID set-up without variation in treatment timing but with two treatments: some groups start receiving the first treatment at a date T^1 , and a subset of those groups then start receiving the second treatment at a later date T^2 . In the absence of the second treatment, one can show that the coefficient on $D_{g,t}^1$ in the regression of $Y_{g,t}$ on group fixed effects, period fixed effects, and $D_{g,t}^1$ identifies the ATT of $D_{g,t}^1$ and does not have negative weights attached to it. On the other hand, in the presence of the second treatment, one can show that β_{fe} no longer identifies the ATT of $D_{g,t}^1$ and may have negative weights attached to it (see Corollary 1 in de Chaisemartin and d'Haultfoeuille, 2021b, a previous version of this paper, for a formal statement and a proof).

The second reason why β_{fe} may differ from δ_{ATT} is that β_{fe} may also be contaminated by the effect of $D_{g,t}^2$: if that effect is heterogeneous across (g, t) cells, $E \left[\sum_{(g,t):D_{g,t}^2=1} (N_{g,t}/N_1) w_{g,t} \Delta_{g,t}^2 \right]$ may differ from zero. Such a contamination phenomenon is not present in the presence of one treatment only (see de Chaisemartin and D'Haultfoeuille, 2020). Below, we give some intuition as to why it arises.

Theorem 1 can be extended to non-binary ordered treatments, that may be continuous or discrete. When $D_{g,t}^1 \neq 0$, let $S_{g,t}^1 = (Y_{g,t}(D_{g,t}^1, D_{g,t}^2) - Y_{g,t}(0, D_{g,t}^2)) / D_{g,t}^1$ be the slope of cell (g, t) 's potential outcome function, when moving its first treatment from 0 to $D_{g,t}^1$, while keeping its second treatment at its observed value. Similarly, when $D_{g,t}^2 \neq 0$, let $S_{g,t}^2 = (Y_{g,t}(0, D_{g,t}^2) - Y_{g,t}(0, 0)) / D_{g,t}^2$.

Finally, let

$$w_{g,t}^k = \frac{\varepsilon_{g,t} D_{g,t}^k}{\sum_{(g,t)} (N_{g,t}/N_1) \varepsilon_{g,t} D_{g,t}^k},$$

for $k = 1, 2$. If $D_{g,t}^1$ and $D_{g,t}^2$ are non-binary, one can show, following similar steps as in the proof of Theorem 1, that

$$\beta_{fe} = E \left[\sum_{(g,t): D_{g,t}^1 \neq 0} \frac{N_{g,t}}{N_1} w_{g,t}^1 S_{g,t}^1 + \sum_{(g,t): D_{g,t}^2 \neq 0} \frac{N_{g,t}}{N_1} w_{g,t}^2 S_{g,t}^2 \right].$$

Moreover, $\sum_{(g,t): D_{g,t}^1 \neq 0} (N_{g,t}/N_1) w_{g,t}^1 = 1$ and $\sum_{(g,t): D_{g,t}^2 \neq 0} (N_{g,t}/N_1) w_{g,t}^2 = 0$. Essentially, Theorem 1 extends to non-binary treatments, replacing the average treatment effects $\Delta_{g,t}^1$ and $\Delta_{g,t}^2$ by slopes of (g, t) -cells' potential outcome functions, from a treatment of zero to their actual treatment. The decomposition in the previous display does not assume a linear treatment effect.

3.1.2 More than two treatment variables

We now go back to the general case where K may be greater than 2. We let $\mathbf{0}^{-1} = (0, \dots, 0)$ be the vector of $K - 1$ zeros. We also define

$$\begin{aligned} \Delta_{g,t}^1 &= Y_{g,t}(1, D_{g,t}^{-1}) - Y_{g,t}(0, D_{g,t}^{-1}), \\ \Delta_{g,t}^{-1} &= Y_{g,t}(0, D_{g,t}^{-1}) - Y_{g,t}(0, \mathbf{0}^{-1}). \end{aligned}$$

$\Delta_{g,t}^1$ is the effect, in cell (g, t) , of moving the first treatment from zero to one while keeping the other treatments at their observed values. $\Delta_{g,t}^{-1}$ is the effect, in cell (g, t) , of moving the other treatments from zero to their actual values, while keeping the first treatment at zero.

Theorem 2 below generalizes Theorem 1.

Theorem 2 *Suppose that Assumptions 1-3 hold. Then,*

$$\beta_{fe} = E \left[\sum_{(g,t): D_{g,t}^1 = 1} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t}^1 + \sum_{(g,t): D_{g,t}^{-1} \neq \mathbf{0}^{-1}} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t}^{-1} \right].$$

Moreover, $\sum_{(g,t): D_{g,t}^1 = 1} (N_{g,t}/N_1) w_{g,t} = 1$, and if $K = 2$ or the treatments $D_{g,t}^2, \dots, D_{g,t}^K$ are mutually exclusive, $\sum_{(g,t): D_{g,t}^{-1} \neq \mathbf{0}^{-1}} (N_{g,t}/N_1) w_{g,t} = 0$.

Theorem 2 is similar to Theorem 1, except that when $K > 2$, we do not always have

$$\sum_{(g,t): D_{g,t}^{-1} \neq \mathbf{0}^{-1}} \frac{N_{g,t}}{N_1} w_{g,t} = 0.$$

The contamination weights on the effects of the other treatments may not sum to 0. Accordingly, even if the effects of all treatments are constant, $\hat{\beta}_{fe}$ may still be biased for the first treatment's effect.

There are three special cases where the weights on the effects of the other treatments sum to 0. The first one is when $K = 2$, as shown in Theorem 1. The second one is when the treatments $D_{g,t}^2, \dots, D_{g,t}^K$ are mutually exclusive, as stated in Theorem 2. The third one is when there is no complementarity or substitutability between the treatments $D_{g,t}^2, \dots, D_{g,t}^K$. Specifically, assume that for all (g, t) , there exists $(\delta_{g,t}^k)_{k=2, \dots, K}$ such that

$$E[\Delta_{g,t}^{-1} | \mathbf{D}] = \sum_{k=2}^K D_{g,t}^k \delta_{g,t}^k. \quad (4)$$

Then, we obtain Decomposition (5) below. The corresponding weights can be computed using the `twowayfeweights` Stata command.

Corollary 1 *Suppose that Assumptions 1-3 and (4) hold. Then,*

$$\beta_{fe} = E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t}^1 + \sum_{k=2}^K \sum_{(g,t):D_{g,t}^k=1} \frac{N_{g,t}}{N_1} w_{g,t} \delta_{g,t}^k \right]. \quad (5)$$

Moreover, $\sum_{(g,t):D_{g,t}^1=1} (N_{g,t}/N_1) w_{g,t} = 1$, and $\sum_{(g,t):D_{g,t}^k=1} (N_{g,t}/N_1) w_{g,t} = 0$ for every $k \in \{2, \dots, K\}$.

On the other hand, when the treatments are not mutually exclusive and may be complementary or substitutable $\hat{\beta}_{fe}$ could be biased even under constant treatment effects. This is because in that case, Regression 1 is misspecified, and should include the interactions of the treatments.

Importantly, Theorem 2 does not necessarily rule out dynamic effects of past treatments on the outcome. The treatments in the regression may for instance be the current treatment and its first $K - 1$ lags. In that case, our potential outcome notation allows the current treatment and its first $K - 1$ lags to affect the outcome. Accordingly, the `twowayfeweights` Stata command can also be used to compute the weights attached to distributed-lags regressions of an outcome on the current treatment and its lags.

3.2 Intuition for, and a perhaps surprising implication of, the contamination bias

3.2.1 Intuition for the contamination bias

The reasons why TWFE regressions are not robust to heterogeneous treatment effects are now well understood (see de Chaisemartin and D'Haultfœuille, 2018; de Chaisemartin and D'Haultfœuille, 2020; Goodman-Bacon, 2021; Borusyak and Jaravel, 2017). In this section, we give intuition as to why β_{fe} may be affected by contamination bias. To do so, we start by considering two very simple examples, one where contamination bias is absent, and the other where it is present.

First, assume that there are three groups and two time periods. With probability one, no group is treated at period 1, and at period 2 group 2 receives the first treatment while group 3 receives the second treatment. Then, it is easy to show that

$$\widehat{\beta}_{fe} = Y_{2,2} - Y_{2,1} - (Y_{1,2} - Y_{1,1}). \quad (6)$$

The right-hand side of the previous display is a DID comparing the period-one-to-two outcome evolution of group 2, that starts receiving the first treatment at period 2, to that of group 1, that is untreated at both dates. Therefore,

$$\begin{aligned} \beta_{fe} &= E(Y_{2,2}(1,0) - Y_{2,1}(0,0) - (Y_{1,2}(0,0) - Y_{1,1}(0,0))) \\ &= E(Y_{2,2}(1,0) - Y_{2,2}(0,0)) + E(Y_{2,2}(0,0) - Y_{2,1}(0,0) - (Y_{1,2}(0,0) - Y_{1,1}(0,0))) \\ &= E(Y_{2,2}(1,0) - Y_{2,2}(0,0)), \end{aligned} \quad (7)$$

where the second equality follows from Assumption 3. Equation (7) is a special case of Equation (3) in Theorem 1. In this simple example, β_{fe} is not contaminated by the effect of the second treatment. It identifies the effect, in group 2 and at period 2, of moving the first treatment from zero to one while keeping the second treatment at its observed value (zero). Because only group 2 at period 2 receives the first treatment, this effect is equal to δ_{ATT} , the ATT of the first treatment controlling for the second treatment.

Now let us consider another example, very similar to that above, but with a fourth group that receives both treatments at period 2. Then, using the equivalence between TWFE regressions and first-difference regressions with two periods and the fact that the first difference of the two treatments are uncorrelated, we obtain

$$\widehat{\beta}_{fe} = \frac{1}{2} (Y_{2,2} - Y_{2,1} - (Y_{1,2} - Y_{1,1})) + \frac{1}{2} (Y_{4,2} - Y_{4,1} - (Y_{3,2} - Y_{3,1})). \quad (8)$$

The first DID in Equation (8) is the same as that in the right-hand side of Equation (6) and it is unbiased for $E(Y_{2,2}(1,0) - Y_{2,2}(0,0))$. The second DID compares the period-one-to-two outcome evolution of group 4, that starts receiving the first and second treatments at period 2, to that of group 3, that only starts receiving the second treatment. Therefore,

$$\begin{aligned} &E(Y_{4,2} - Y_{4,1} - (Y_{3,2} - Y_{3,1})) \\ &= E(Y_{4,2}(1,1) - Y_{4,1}(0,0) - (Y_{3,2}(0,1) - Y_{3,1}(0,0))) \\ &= E(Y_{4,2}(1,1) - Y_{4,2}(0,1)) + E(Y_{4,2}(0,1) - Y_{4,2}(0,0)) - E(Y_{3,2}(0,1) - Y_{3,2}(0,0)) \\ &+ E(Y_{4,2}(0,0) - Y_{4,1}(0,0) - (Y_{3,2}(0,0) - Y_{3,1}(0,0))) \\ &= E(Y_{4,2}(1,1) - Y_{4,2}(0,1)) + E(Y_{4,2}(0,1) - Y_{4,2}(0,0)) - E(Y_{3,2}(0,1) - Y_{3,2}(0,0)). \end{aligned} \quad (9)$$

Equations (8) and (9) imply that

$$\begin{aligned} \beta_{fe} = & \frac{1}{2}E(Y_{2,2}(1,0) - Y_{2,2}(0,0)) + \frac{1}{2}E(Y_{4,2}(1,1) - Y_{4,2}(0,1)) \\ & + \frac{1}{2}E(Y_{4,2}(0,1) - Y_{4,2}(0,0)) - \frac{1}{2}E(Y_{3,2}(0,1) - Y_{3,2}(0,0)). \end{aligned} \quad (10)$$

Equation (10) is a special case of Equation (3) in Theorem 1. β_{fe} identifies the sum of two terms. The term on the first line is the average effect, in groups two and four and at period two, of moving the first treatment from zero to one while keeping the second treatment at its observed value (zero in group 2, one in group 4). The term on the second line is a contamination bias term, equal to the difference, between groups 4 and 3, of the effect of moving the second treatment from zero to one while keeping the first treatment at zero.

The contamination bias appears in the second example because $\widehat{\beta}_{fe}$ leverages a DID comparing a group that starts receiving the first and the second treatments to a group that starts receiving the second treatment only. With heterogeneous treatment effects, this comparison is contaminated by the effect of the second treatment. On the other hand, if the effect of the second treatment does not vary across groups, this contamination bias disappears. To our knowledge, our paper is the first to show that TWFE regressions with several treatments leverage this type of “forbidden comparisons”, using the terminology coined by Borusyak and Jaravel (2017).

In the example with four groups, a simple solution to eliminate the contamination bias is to add the interaction of the two treatments to the regression. One can in fact show the following, slightly more general result. With only two time periods, and groups that do not receive any of the two treatments in the first period, the coefficient on $D_{g,t}^1$ in the regression of $Y_{g,t}$ on $D_{g,t}^1$, $D_{g,t}^2$, and $D_{g,t}^1 D_{g,t}^2$ is not contaminated by the effect of the second treatment. In such cases, the regression with the interaction term is preferable, as it makes the contamination problem disappear. This result does not, however, translate to more general designs with more than two time periods and where groups may receive the treatments at every period. It is easy to find examples where adding the interaction to the regression actually increases the contamination weights. This is the case for instance in the application we consider in Section 5: in the regression without control variables and with the two main treatments (the minimum staff-to-child ratio and the minimum number of years of schooling required for daycare directors), adding the interaction between the two treatments actually increases the absolute value of the contamination weights.

In the first example, the two treatments are mutually exclusive so $\widehat{\beta}_{fe}$ cannot leverage a “forbidden” DID comparing a group that starts receiving the first and the second treatments to a group that starts receiving the second treatment only, which is why there is no contamination bias in this example. This does not mean contamination bias never arises with mutually exclusive treatments. To illustrate this point, let us consider a third example with two groups and three periods. Group 1 receives the first treatment at period 3, and Group 2 receives the

second treatment at periods 2 and 3. Then, because this regression is equivalent to a regression of $Y_{2,t} - Y_{1,t}$ on a constant, $D_{2,t}^1 - D_{1,t}^1$ and $D_{2,t}^2 - D_{1,t}^2$, we obtain, after some algebra,

$$\widehat{\beta}_{fe} = Y_{1,3} - Y_{1,2} - (Y_{2,3} - Y_{2,2}). \quad (11)$$

Accordingly, one can show that

$$\begin{aligned} \beta_{fe} = & E(Y_{1,3}(1, 0) - Y_{1,3}(0, 0)) \\ & + E(Y_{2,2}(0, 1) - Y_{2,2}(0, 0)) - E(Y_{2,3}(0, 1) - Y_{2,3}(0, 0)). \end{aligned} \quad (12)$$

$\widehat{\beta}_{fe}$ is contaminated by the effect of the second treatment, because it leverages a DID where the control group receives the second treatment at both dates. This second type of “forbidden” DID is very similar to the late- versus early-treated DIDs due to which TWFE regressions with one treatment are not robust to heterogeneous treatment effects (see de Chaisemartin and D’Haultfœuille, 2020; Goodman-Bacon, 2021; Borusyak and Jaravel, 2017). Note that if the effect of the second treatment is constant over time, the contamination bias term disappears. However, constant effects over time is often an implausible assumption.

Overall, TWFE regressions with several treatments are not affected by contamination bias in very simple designs with two time periods, where groups are only treated in the second period, and where the treatments are mutually exclusive. In designs with non-mutually exclusive treatments, contamination bias may appear because $\widehat{\beta}_{fe}$ may leverage DIDs comparing a group that starts receiving, say, the first and the second treatments to a group that starts receiving the second treatment only. With more than two time periods, even if the treatments are mutually exclusive, $\widehat{\beta}_{fe}$ may leverage DIDs comparing a group that starts receiving, say, the first treatment, to a group receiving the second treatment at both dates.

3.2.2 A perhaps surprising implication of the contamination bias

Theorem 1 has an important and perhaps surprising consequence for TWFE regressions with one treatment where one seeks to estimate heterogeneous treatment effects. Oftentimes, researchers run a TWFE regression with a treatment variable $D_{g,t}$ interacted with a group-level binary variable I_g , and with $(1 - I_g)$.³ For instance, to study if the treatment effect differs in poor and rich counties, one interacts the treatment with an indicator for counties above the median income, and with an indicator for counties below the median income. Theorem 1 also applies to those regressions. Specifically, one has

$$\beta_{fe}^{I=1} = E \left[\sum_{(g,t):D_{g,t}=1,I_g=1} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t} + \sum_{(g,t):D_{g,t}=1,I_g=0} \frac{N_{g,t}}{N_1} w_{g,t} \Delta_{g,t} \right].$$

³Researchers may instead have $D_{g,t}$ and $D_{g,t}I_g$ in the regression. The coefficient on $D_{g,t}$ in this regression is equal to that on $D_{g,t}(1 - I_g)$ in the regression described in the text. The coefficient on $D_{g,t}I_g$ is equal to the difference between that on $D_{g,t}I_g$ and that on $D_{g,t}(1 - I_g)$ in the regression described in the text. Accordingly, the discussion in this section also applies to those regressions.

where $\beta_{fe}^{I=1}$ is the coefficient on $D_{g,t} \times I_g$, and $\Delta_{g,t} = Y_{g,t}(1) - Y_{g,t}(0)$. The previous display implies that the coefficient on $D_{g,t} \times I_g$ is contaminated by the treatment effect in (g, t) cells such that $I_g = 0$. In the example, the coefficient on the treatment interacted with the indicator for rich counties is contaminated by the treatment effect in poor counties. This calls into question the use of such TWFE regressions to estimate heterogeneous effects.

This contamination phenomenon disappears if the time fixed effects are interacted with I_g in the regression. Then, the coefficient on $D_{g,t} \times I_g$ becomes equivalent to that one would obtain by running a TWFE regression restricting the sample to groups such that $I_g = 1$. It follows from de Chaisemartin and D'Haultfoeulle (2020) that this coefficient identifies a weighted sum of the treatment effects across (g, t) cells such that $D_{g,t} = 1, I_g = 1$: it is not contaminated by the treatment effect in (g, t) cells such that $D_{g,t} = 1, I_g = 0$.

3.3 Should one control for other treatments?

In this section, we derive a decomposition similar to that in Theorem 1, when there are two treatments but the second treatment is omitted from the regression.

Regression 2 (Short TWFE regression)

Let $\beta_{fe}^s = E[\widehat{\beta}_{fe}^s]$, where $\widehat{\beta}_{fe}^s$ denotes the coefficient on $D_{g,t}^1$ in a sample OLS regression of $Y_{g,t}$ on group fixed effects, period fixed effects, and $D_{g,t}^1$, weighted by $N_{g,t}$.

Let $\varepsilon_{g,t}^s$ denote the residual of cell (g, t) in the sample regression of $D_{g,t}^1$ on group and period fixed effects. If the regressors in Regression 2 are not collinear, the average value of $\varepsilon_{g,t}^s$ across all (g, t) cells with $D_{g,t}^1 = 1$ differs from 0: $\sum_{(g,t):D_{g,t}^1=1} (N_{g,t}/N_1) \varepsilon_{g,t}^s \neq 0$. Then we let $w_{g,t}^s$ denote $\varepsilon_{g,t}^s$ divided by that average:

$$w_{g,t}^s = \frac{\varepsilon_{g,t}^s}{\sum_{(g,t):D_{g,t}^1=1} (N_{g,t}/N_1) \varepsilon_{g,t}^s}.$$

Theorem 3 *Suppose that Assumptions 1-3 hold and $K = 2$. Then,*

$$\beta_{fe}^s = E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} w_{g,t}^s \Delta_{g,t}^1 + \sum_{(g,t):D_{g,t}^2=1} \frac{N_{g,t}}{N_1} w_{g,t}^s \Delta_{g,t}^2 \right]. \quad (13)$$

Moreover, $\sum_{(g,t):D_{g,t}^1=1} (N_{g,t}/N_1) w_{g,t}^s = 1$ and $\sum_{(g,t):D_{g,t}^2=1} (N_{g,t}/N_1) w_{g,t}^s$ may differ from zero.

Theorem 3 is similar to Theorem 1. It shows that the coefficient on $D_{g,t}^1$ in the short regression identifies the sum of two terms. The first term in Theorem 3 is similar to that in Theorem 1, namely a weighted sum of the effect of moving $D_{g,t}^1$ from 0 to 1 while keeping $D_{g,t}^2$ at its observed value, but with different weights that still sum to one. The second term in Theorem 3 is also

similar to that in Theorem 1, namely a weighted sum of the effect of moving $D_{g,t}^2$ from 0 to 1 while keeping $D_{g,t}^1$ at 0, but with different weights that no longer sum to zero. Note that Theorem 3 can easily be extended to instances with more than two treatments.

Theorem 3 may be seen as a version of the standard omitted variable bias (OVB) formula, for TWFE regressions and with heterogeneous treatment effects. Allowing for heterogeneous treatment effects in this OVB formula has one important consequence. With constant treatment effects, there is no OVB if the omitted variable $D_{g,t}^2$ is uncorrelated with $\varepsilon_{g,t}^s$, the residual of $D_{g,t}^1$ from a regression on group and time FEs. On the other hand, with heterogeneous treatment effects, there may be an OVB even if $D_{g,t}^2$ and $\varepsilon_{g,t}^s$ are uncorrelated, if $\sum_{(g,t)}(N_{g,t}/N_1)w_{g,t}^s D_{g,t}^2 \Delta_{g,t}^2 \neq 0$, namely if $D_{g,t}^2 \Delta_{g,t}^2$ and $\varepsilon_{g,t}^s$ are correlated. Conversely, there is no OVB if $D_{g,t}^2 \Delta_{g,t}^2$ and $\varepsilon_{g,t}^s$ are uncorrelated, but this condition is strong, and unlike the standard condition that $D_{g,t}^2$ and $\varepsilon_{g,t}^s$ be uncorrelated, it is untestable. The implication of this result is that in TWFE regressions, if treatment effects are heterogeneous, failing to control for other treatments may lead to an OVB even if those other treatments are uncorrelated with the main treatment of interest.⁴

The previous paragraph may seem to imply that in TWFE regressions, one should control for all the time-varying treatments one observes, if not for all that exist. And indeed, under constant effects and a parallel trends assumption on $Y_{g,t}(0,0)$, omitting the second treatment from the regression leads to an omitted variable bias, and including the second treatment into the regression is always preferable. But again, this logic breaks down under heterogeneous treatment effects. Then, including the second treatment into the regression may not be preferable: $D_{g,t}^1$'s coefficient in the long regression may be more biased for δ_{ATT} than $D_{g,t}^1$'s coefficient in the short regression. The following corollary formalizes this idea.

Corollary 2 *Suppose that Assumptions 1-3 hold, $K = 2$, and there is a real number B such that $|\Delta_{g,t}^1| \leq B$ and $|\Delta_{g,t}^2| \leq B$ for all (g, t) . Then,*

$$|\beta_{fe} - \delta_{ATT}| \leq B \times E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} |w_{g,t} - 1| + \sum_{(g,t):D_{g,t}^2=1} \frac{N_{g,t}}{N_1} |w_{g,t}| \right],$$

$$|\beta_{fe}^s - \delta_{ATT}| \leq B \times E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} |w_{g,t}^s - 1| + \sum_{(g,t):D_{g,t}^2=1} \frac{N_{g,t}}{N_1} |w_{g,t}^s| \right].$$

Moreover, both upper bounds are sharp.

⁴Not all OLS regressions are subject to this issue in the presence of heterogeneous treatment effects. For instance, in an RCT, if one regresses the outcome on the treatment omitting other determinants of the outcome, there will be no OVB in the regression, because the effect of those other determinants of the outcome is by design uncorrelated with the treatment assignment. In TWFE regressions on the other hand, there is no guarantee that $\sum_{(g,t)}(N_{g,t}/N_1)w_{g,t}^s D_{g,t}^2 \Delta_{g,t}^2 = 0$.

Corollary 2 assumes that the effects of the first and second treatments are both bounded in every (g, t) cell by a constant B . Under that assumption, it gives the maximal biases of $\widehat{\beta}_{fe}$ and $\widehat{\beta}_{fe}^s$ as estimators of δ_{ATT} , the ATT of $D_{g,t}^1$ controlling for $D_{g,t}^2$. One can compare those maximal biases by comparing (estimates of)

$$E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} |w_{g,t} - 1| + \sum_{(g,t):D_{g,t}^2=1} \frac{N_{g,t}}{N_1} |w_{g,t}| \right]$$

and

$$E \left[\sum_{(g,t):D_{g,t}^1=1} \frac{N_{g,t}}{N_1} |w_{g,t}^s - 1| + \sum_{(g,t):D_{g,t}^2=1} \frac{N_{g,t}}{N_1} |w_{g,t}^s| \right],$$

which does not require specifying B .⁵ The maximal bias of $\widehat{\beta}_{fe}$ could be larger than that of $\widehat{\beta}_{fe}^s$, if for (g, t) s such that $D_{g,t}^1 = 1$ the weights $w_{g,t}$ are on average further away from one than the weights $w_{g,t}^s$, and/or if for (g, t) s such that $D_{g,t}^2 = 1$ the contamination weights $w_{g,t}$ are on average further away from zero than the weights $w_{g,t}^s$. In our application in Section 5, we find that the estimated maximal bias of the long regression is almost five times larger than that of the short regression. Then, the short regression is preferable, at least per our maximal-bias metric.

3.4 Related literature

3.4.1 TWFE regressions with control variables

Theorem S4 in the Web Appendix of de Chaisemartin and D'Haultfœuille (2020) studies TWFE regressions with one treatment and some time-varying control variables $X_{g,t}$.⁶ It assumes that for some vector θ , $E(Y_{g,t}(0) - Y_{g,t-1}(0) - (X_{g,t} - X_{g,t-1})'\theta)$ does not vary across groups. In other words, groups may experience different trends, but those are fully accounted for by a linear model in the evolution of their covariates. Under that assumption, de Chaisemartin and D'Haultfœuille (2020) show that TWFE regressions with one treatment and some controls identify a weighted sum of the treatment effects across all treated (g, t) cells, with weights that differ from those without controls.

Theorems 1 and 2 are related to, but different from, that result. For instance, with $D_{g,t}^2$ as the only control variable in the regression, the weighted sum of treatment effects identified by β_{fe} in Theorem S4 of de Chaisemartin and D'Haultfœuille (2020) is identical to the first weighted sum in Theorem 1. On the other hand, the second weighted sum in Theorem 1, the

⁵A similar result holds if we consider distinct bounds B_1 and B_2 for $|\Delta_{g,t}^1|$ and $|\Delta_{g,t}^2|$. Then, one has to multiply $\sum_{(g,t):D_{g,t}^2=1} (N_{g,t}/N_1) |w_{g,t}|$ by (B_2/B_1) when performing the comparison of the maximal biases. Hence, in this case, one needs to take a stand on the ratio B_2/B_1 .

⁶With time-invariant group-level controls, the TWFE regression is not identified.

contamination term, does not appear in Theorem S4 of de Chaisemartin and D’Haultfœuille (2020). The only case where the decompositions in Theorems 1 and 2 reduce to that in Theorem S4 of de Chaisemartin and D’Haultfœuille (2020) is when the effect of the other treatments do not vary across (g, t) cells, an often implausible assumption.

Treating a variable in Regression 1 as a covariate or as a treatment leads to different decompositions of β_{fe} , so it is important to weigh this choice carefully. In principle, any variable that has a causal effect on the outcome, in the sense that different values of that variable lead to different counterfactual outcomes, should be regarded as a treatment. Covariates, on the other hand, are variables that do not have a causal effect on the outcome, though their evolution may be correlated with the outcome’s evolution.

Overall, Theorems 1 and 2 are important extensions of Theorem S4 of de Chaisemartin and D’Haultfœuille (2020). They apply to TWFE regressions where some of the variables in the regression, other than the main treatment of interest, can affect the outcome and may have heterogeneous effects. Such regressions are likely to be common in practice.

3.4.2 *Linear regressions with several treatments*

Theorems 1 and 2 complement the pioneering work of Sun and Abraham (2021). The authors study the so-called event-study regression, an example of a TWFE regression with several treatments, where the treatments are indicators for having started receiving a single binary-and-staggered treatment ℓ periods ago. In those regressions, the authors show that effects of being treated for ℓ' periods may contaminate the coefficient supposed to measure the effect of ℓ periods of treatment in the regression, and they provide a decomposition formula one can use to quantify the extent of the phenomenon. If i) the K treatments in Regression 1 are indicators for having started receiving a single binary-and-staggered treatment ℓ periods ago, and ii) the treatment no longer has an effect after $K + 1$ periods of exposure, then our Theorem 2 reduces to Proposition 3 in Sun and Abraham (2021), provided no lags are gathered together in the event-study regression they consider.⁷ Our decompositions extend their result, by showing that the contamination bias they first uncovered is much more pervasive: it can arise in any TWFE regression with several treatments, rather than in event-study regressions only. In particular, our results apply to situations where the treatments are different, potentially non-mutually exclusive policies, that may not be binary or may not follow a staggered adoption design. Another difference with their work is that with non-mutually exclusive treatments, the contamination weights do not sum to zero. We also provide some novel intuition as to why contamination may arise with different, potentially non-mutually exclusive treatments in the regression. We also

⁷In their decomposition, Sun and Abraham (2021) gather groups that started receiving the treatment at the same period into cohorts. Their decomposition can then be further decomposed, finally leading to the result in our Theorem 2.

show that TWFE regressions with one treatment intended at measuring heterogeneous treatment effects may also be affected by some form of contamination bias. Finally, we show that omitting the other treatments from the regression may not necessarily increase the regression coefficient’s bias.

Theorem 1 is also related to the pioneering work of Hull (2018). In his Section 2.2, the author studies TWFE regressions where indicators for each value that a multinomial treatment may take are included in the regression, an example of a TWFE regression with several treatments. Equation (15) therein is, to our knowledge, the first instance where a contamination phenomenon was shown. However, the paper does not discuss this phenomenon. It also does not give a decomposition formula like Theorem 1, so one cannot use the paper’s results to compute the contamination weights, and assess whether they are important in a given regression. Finally, the paper’s result applies when the data has two periods, and in instances where the treatments in the regression are indicators for each value that a multinomial treatment may take.

Another related paper, released after ours, is Goldsmith-Pinkham, Hull and Kolesár (2021), who show that a contamination phenomenon similar to that in Sun and Abraham (2021) and in Theorem 1 also arises in linear regressions with several treatments, and a set of controls such that the treatments can be assumed to be independent of the potential outcomes conditional on those controls. Their result is not nested within and does not nest the results of Sun and Abraham (2021) nor ours: both Sun and Abraham (2021) and us assume parallel trends rather than conditional independence. The weights in their decomposition are functions of the variance-covariance matrix of the treatments conditional on the controls. An interesting difference with our results is that under their conditional independence assumption, the weights on the effect of $D_{g,t}^1$ are all positive.

Overall, our four papers complement each other, and show that the contamination phenomenon is very pervasive, as it arises under several identifying assumptions (parallel trends and conditional independence), and irrespective of the nature of the treatments included in the regression.

4 Alternative estimator

4.1 Identifying assumption

In this section, we start by considering the following identifying assumption.

Assumption 4 (*Strong exogeneity and common trends from $t - 1$ to t , conditional on $D_{g,t-1}$*)
For all $(g, t) \in \{1, \dots, G\} \times \{2, \dots, T\}$ and all $d_{t-1} \in \{0, 1\}^K$,

1. $E(Y_{g,t}(d_{t-1}) - Y_{g,t-1}(d_{t-1}) | D_{g,1}, \dots, D_{g,t-2}, D_{g,t-1} = d_{t-1}, D_{g,t}, \dots, D_{g,T}) = E(Y_{g,t}(d_{t-1}) - Y_{g,t-1}(d_{t-1}) | D_{g,t-1} = d_{t-1})$.

2. $E(Y_{g,t}(d_{t-1}) - Y_{g,t-1}(d_{t-1}) | D_{g,t-1} = d_{t-1})$ does not vary across g .

Like Assumption 3, Assumption 4 imposes both a strong exogeneity and a parallel trends condition. The strong exogeneity condition requires that groups' $t - 1$ -to- t outcome evolution, in the counterfactual scenario where their period- t treatments all remain at their $t - 1$ value, be mean independent of their treatments at every period other than $t - 1$. The parallel trends assumption requires that groups with the same period- $t - 1$ treatments have the same counterfactual trends. Then, consider a group whose first treatment changes between $t - 1$ and t , but whose other treatments remain constant. Under Assumption 4, the $t - 1$ -to- t evolution of its outcome had its first treatment not changed is identified by the outcome evolution of groups whose treatments all remain constant and with the same period- $t - 1$ treatments.

We now compare our new assumption, Assumption 4, to the more standard Assumption 3. The two assumptions are non-nested, and there are two main differences between them. First, Assumption 3 requires that all groups be on parallel trends, over the entire duration of the panel. Assumption 4, on the other hand, only requires that groups with the same period- $t - 1$ treatments be on parallel trends, from $t - 1$ to t . Assumption 4 may then be more plausible: groups with the same treatments in the baseline period may be more similar, and may be more likely to experience parallel trends.⁸ Moreover, parallel trends may be more likely to hold over consecutive time periods than over the panel's entire duration.

Second, Assumption 3 is a parallel trends assumption in the counterfactual where groups do not receive any treatment, while Assumption 4 is a parallel trends assumption in the counterfactual where groups' treatments do not change from $t - 1$ to t . Accordingly, Assumption 3 only restricts one potential outcome, the one without any treatment, while Assumption 4 imposes restrictions on many potential outcomes. Still, Assumption 4 does not impose any restriction on treatment effect heterogeneity, because it restricts only one potential outcome per (g, t) cell, namely $Y_{g,t}(d_{t-1})$ for (g, t) cells such that $D_{g,t-1} = d_{t-1}$. In particular, Assumption 4 does not require that all groups experience the same evolution of their treatment effect. Moreover, in complicated designs where the number of treatments is large and/or when the treatments are non binary, Assumption 4 may have considerably more identifying power than Assumption 3. Under Assumption 3, an heterogeneity-robust DID estimator can only use as controls groups that do not receive any treatment at two dates at least. Moreover, treatment effects can only be estimated

⁸Because it imposes parallel trends conditional on $D_{g,t-1}$, Assumption 4 may be seen as “in-between” a standard parallel trends assumption and the sequential ignorability assumption, another commonly-used identifying assumption in panel data models (see, e.g., Robins, 1986; Bojinov, Rambachan and Shephard, 2021). Sequential ignorability requires that treatment be uncounfounded conditional on prior treatment and outcome, which implies parallel trends conditional on prior treatment and outcome. Because Assumption 4 does not condition on groups' $t - 1$ outcomes, it may be less plausible than sequential ignorability. At the same time, estimators relying on sequential ignorability need to compare groups with the same prior treatments and outcomes. This may lead to a curse of dimensionality.

for groups that do not receive any treatment at one date at least. With many treatments and/or when the treatments are non binary, those two sets of groups may be small. In our empirical application in Section 5, there are two non-binary treatments, and while there are (g, t) cells whose two treatments are equal to 0, there is no group that does not receive any of the two treatments at two dates at least. Accordingly, we cannot construct an heterogeneity-robust DID estimator relying on Assumption 3, while we can construct one relying on Assumption 4.

We also consider a second identifying assumption.

Assumption 5 (*Strong exogeneity and common trends from $t - 1$ to t , conditional on $D_{g,t}$*) For all $(g, t) \in \{1, \dots, G\} \times \{2, \dots, T\}$ and all $d_t \in \{0, 1\}^K$,

1. $E(Y_{g,t}(d_t) - Y_{g,t-1}(d_t) | D_{g,1}, \dots, D_{g,t-1}, D_{g,t} = d_t, D_{g,t+1}, \dots, D_{g,T}) = E(Y_{g,t}(d_t) - Y_{g,t-1}(d_t) | D_{g,t} = d_t)$.
2. $E(Y_{g,t}(d_t) - Y_{g,t-1}(d_t) | D_{g,t} = d_t)$ does not vary across g .

Assumption 5 is similar to Assumption 4, except that it assumes parallel trends from $t - 1$ to t , in the counterfactual where groups keep their period- t rather than their period- $t - 1$ treatments. Assumptions 4 and 5 are both parallel trends assumptions over two consecutive periods, among groups with the same treatments at one of the two periods. Accordingly, in instances where Assumption 4 is plausible, Assumption 5 may be plausible too. Imposing jointly Assumptions 4 and 5 may imply that the treatment effects follow the same evolution over time in some groups.⁹

4.2 Target parameters

Let

$$\mathcal{S}_1 = \left\{ (g, t) : t \geq 2, D_{g,t}^1 \neq D_{g,t-1}^1, D_{g,t}^{-1} = D_{g,t-1}^{-1}, \exists g' : D_{g',t} = D_{g',t-1} = D_{g,t-1} \right\}$$

and $N_{\mathcal{S}_1} = \sum_{(g,t) \in \mathcal{S}_1} N_{g,t}$. \mathcal{S}_1 is the set of cells (g, t) whose first treatment changes between $t - 1$ and t while their other treatments do not change, and such that there is another group g' whose treatments do not change between $t - 1$ and t , and with the same treatments as g in $t - 1$. Hereafter, those cells are referred to as switchers. We show below that under Assumption 4, one can unbiasedly estimate

$$\delta_1 = E \left[\sum_{(g,t) \in \mathcal{S}_1} \frac{N_{g,t}}{N_{\mathcal{S}_1}} \Delta_{g,t}^1 \right],$$

⁹For instance, if $K = 1$, $G = 4$, $T = 2$, $D_{1,1}^1 = D_{1,2}^1 = 0$, $D_{2,1}^1 = D_{2,2}^1 = 1$, $D_{3,1}^1 = 0, D_{3,2}^1 = 1$, and $D_{4,1}^1 = 1, D_{4,2}^1 = 0$, one can show that together, Assumptions 4 and 5 imply that the treatment effect follows the same evolution in groups 3 and 4.

the average effect of moving the first treatment from 0 to 1 while keeping all other treatments at their observed value, across all switchers.¹⁰

δ_1 may differ from δ_{ATT} , arguably a more natural target parameter. The two parameters apply to different and non-nested sets of (g, t) cells. Let $\mathcal{D}_1 = \{(g, t) : D_{g,t}^1 = 1\}$. δ_1 is the average of $\Delta_{g,t}^1$ across all cells in \mathcal{S}_1 . δ_{ATT} is the average effect of $\Delta_{g,t}^1$ across all cells in \mathcal{D}_1 .

(g, t) cells belonging to \mathcal{D}_1 but not to \mathcal{S}_1 can be divided into five mutually exclusive subgroups, detailed in Section 1 of the Web Appendix. Identifying the effect of the first treatment in each of those subgroups would require imposing stronger assumptions than Assumption 4. For instance, the first subgroup belonging to \mathcal{D}_1 but not to \mathcal{S}_1 are all (g, t) s such that $D_{g,t}^1 = 1$ for all t . As those cells' first treatment never changes, their first-treatment's effect cannot be identified under a parallel trends assumption. The second and third subgroups are cells whose first-treatment's effect could only be identified under a stronger parallel trends assumption than Assumption 4, which only imposes parallel trends over consecutive periods and conditional on cells' period- $t - 1$ treatments. The fourth subgroup are cells whose first treatment changes while at least one of their other treatment also changes. Identifying their first-treatment's effect would require assuming that the effect of the other treatments is constant between groups, as discussed in Section 3.2.1 (see Equation (10) therein). The fifth subgroup are cells whose first treatment changes while their other treatments do not change, but such that all potential control cells experiencing no treatment change have different baseline treatments. Identifying their first-treatment's effect would require assuming that the effects of the other treatments are constant over time, as discussed in Section 3.2.1 (see Equation (12) therein). Therefore, \mathcal{S}_1 is the maximal set of (g, t) cells for which the effect of the first treatment can be identified under a minimal parallel trends assumption and without restricting treatment effect heterogeneity.

Finally, while we expect \mathcal{S}_1 to be often smaller than \mathcal{D}_1 , there are also (g, t) cells that belong to \mathcal{S}_1 but not to \mathcal{D}_1 . Those are the switching-out cells, such that $D_{g,t}^1 = 0, D_{g,t-1}^1 = 1, D_{g,t}^{-1} = D_{g,t-1}^{-1}, \exists g' : D_{g',t} = D_{g',t-1} = D_{g,t-1}$.

As δ_1 and δ_{ATT} apply to different, non-nested subpopulations, a significant difference between $\widehat{\beta}_{fe}$ and the estimator of δ_1 we propose below cannot be interpreted as evidence that $\widehat{\beta}_{fe}$ is biased for δ_{ATT} . It could also be the case that $\widehat{\beta}_{fe}$ is unbiased for δ_{ATT} and δ_1 and δ_{ATT} differ. On the other hand, under Assumptions 3 and 4, a significant difference between $\widehat{\beta}_{fe}$ and the estimator of δ_1 implies that the effect of at least one treatment is not constant.

Similarly, we show below that under Assumption 5, one can unbiasedly estimate

$$\delta_2 = E \left[\sum_{(g,t) \in \mathcal{S}_2} \frac{N_{g,t}}{N_{\mathcal{S}_2}} \Delta_{g,t}^1 \right],$$

¹⁰When $N_{\mathcal{S}_1} = 0$, we simply let the term inside brackets be equal to 0.

where

$$\mathcal{S}_2 = \left\{ (g, t) : t \leq T - 1, D_{g,t}^1 \neq D_{g,t+1}^1, D_{g,t}^{-1} = D_{g,t+1}^{-1}, \exists g' : D_{g',t} = D_{g',t+1} = D_{g,t+1} \right\},$$

and $N_{\mathcal{S}_2} = \sum_{(g,t) \in \mathcal{S}_2} N_{g,t}$. \mathcal{S}_2 is the set of cells (g, t) whose first treatment changes between t and $t + 1$ while their other treatments do not change, and such that there is another group g' whose treatments do not change between t and $t + 1$, and with the same treatments as g in $t + 1$. \mathcal{S}_1 and \mathcal{S}_2 are not necessarily disjoint: a (g, t) cell experiencing two consecutive changes of its first treatment ($D_{g,t-1}^1 \neq D_{g,t}^1$ and $D_{g,t}^1 \neq D_{g,t+1}^1$) may belong both to δ_1 and to δ_2 . On the other hand, a (g, t) cell that does not experience two consecutive changes of its first treatment ($D_{g,t-1}^1 = D_{g,t}^1$ or $D_{g,t}^1 = D_{g,t+1}^1$) may belong to δ_1 or to δ_2 but cannot belong to both sets.

Finally, under Assumptions 4 and 5, one can unbiasedly estimate

$$\delta = E \left[\sum_{(g,t) \in \mathcal{S}_1 \cup \mathcal{S}_2} \frac{N_{g,t}}{N_{\mathcal{S}_1 \cup \mathcal{S}_2}} \Delta_{g,t}^1 \right],$$

where $N_{\mathcal{S}_1 \cup \mathcal{S}_2} = \sum_{(g,t) \in \mathcal{S}_1 \cup \mathcal{S}_2} N_{g,t}$.

4.3 Estimation

We now show that under Assumption 4, δ_1 can be unbiasedly estimated by a weighted average of DID. For all $t \in \{2, \dots, T\}$, for all $(d, d') \in (\mathcal{D}_1)^2$, and for all $d^{-1} \in \mathcal{D}_2 \times \dots \times \mathcal{D}_K$, let

$$\mathcal{G}_{d,d',d^{-1},t} = \left\{ g : D_{g,t}^1 = d, D_{g,t-1}^1 = d', D_{g,t}^{-1} = D_{g,t-1}^{-1} = d^{-1} \right\}$$

be the set of groups whose first treatment goes from d' to d from $t - 1$ to t while their other treatments are equal to d^{-1} at both dates. We then let $N_{d,d',d^{-1},t} = \sum_{g \in \mathcal{G}_{d,d',d^{-1},t}} N_{g,t}$ denote the total population of groups in $\mathcal{G}_{d,d',d^{-1},t}$. Let also

$$\text{DID}_{+,d^{-1},t}^f = \sum_{g \in \mathcal{G}_{1,0,d^{-1},t}} \frac{N_{g,t}}{N_{1,0,d^{-1},t}} (Y_{g,t} - Y_{g,t-1}) - \sum_{g \in \mathcal{G}_{0,0,d^{-1},t}} \frac{N_{g,t}}{N_{0,0,d^{-1},t}} (Y_{g,t} - Y_{g,t-1}), \quad (14)$$

$$\text{DID}_{-,d^{-1},t}^f = \sum_{g \in \mathcal{G}_{1,1,d^{-1},t}} \frac{N_{g,t}}{N_{1,1,d^{-1},t}} (Y_{g,t} - Y_{g,t-1}) - \sum_{g \in \mathcal{G}_{0,1,d^{-1},t}} \frac{N_{g,t}}{N_{0,1,d^{-1},t}} (Y_{g,t} - Y_{g,t-1}). \quad (15)$$

Note that $\text{DID}_{+,d^{-1},t}^f$ is not defined when $N_{1,0,d^{-1},t} = 0$ or $N_{0,0,d^{-1},t} = 0$. In such instances, we let $\text{DID}_{+,d^{-1},t}^f = 0$. Similarly, we let $\text{DID}_{-,d^{-1},t}^f = 0$ when $N_{1,1,d^{-1},t} = 0$ or $N_{0,1,d^{-1},t} = 0$.

$\text{DID}_{+,d^{-1},t}^f$ compares the $t - 1$ -to- t outcome evolution of groups whose first treatment goes from 0 to 1 from $t - 1$ to t while their other treatments are equal to d^{-1} at both dates, to the outcome evolution of groups whose first and other treatments are respectively equal to 0 and d^{-1} at both dates. Under Assumption 4, the latter evolution is a valid counterfactual of the outcome

evolution that the first groups would have experienced if their first treatment had remained equal to 0 at period t . $\text{DID}_{-,d^{-1},t}^f$'s interpretation is similar, except that it compares groups whose first treatment is equal to 1 at both dates to groups whose first treatment goes from 1 to 0.

Finally, let

$$\text{DID}_M^f = \sum_{t=2}^T \sum_{d^{-1} \in \{0,1\}^{K-1}} \left(\frac{N_{1,0,d^{-1},t}}{N_{S_1}} \text{DID}_{+,d^{-1},t}^f + \frac{N_{0,1,d^{-1},t}}{N_{S_1}} \text{DID}_{-,d^{-1},t}^f \right) \quad (16)$$

if $N_{S_1} > 0$, and $\text{DID}_M^f = 0$ if $N_{S_1} = 0$. DID_M^f is just a weighted average of the $\text{DID}_{+,d^{-1},t}^f$ and $\text{DID}_{-,d^{-1},t}^f$ estimators, across values of the other treatments d^{-1} and across time periods t .

Theorem 4 *If Assumptions 1-2 and 4 hold, $E[\text{DID}_M^f] = \delta_1$.*

DID_M^f extends the DID_M estimator in de Chaisemartin and D'Haultfoeuille (2020) to settings with several treatments. With several treatments, one could show the analogue of Theorem 3 for the DID_M estimator in de Chaisemartin and D'Haultfoeuille (2020): the fact that this estimator does not control for the other treatments may lead to a bias, even if switchers and non-switchers are equally likely to experience a change in their other treatments, the analogue of having that the treatments are uncorrelated conditional on the group and time fixed effects in the TWFE regression. To avoid that, the DID_M^f and DID_M estimators differ on three important dimensions: DID_M^f does not estimate the effect of the first treatment in (g, t) cells such that at least one of g 's other treatments changes between $t - 1$ and t ; it drops control groups whose first treatment does not change but such that at least one of their other treatments changes between $t - 1$ and t ; and it compares switchers and non-switchers with the same baseline values of their other treatments. All those modifications ensure that our new estimator is not biased in the presence of other treatments with potentially heterogeneous treatment effects, but they may also come at a cost in terms of precision: the DID_M^f estimator in this paper discards several cells from the estimation. Accordingly, there may be a bias-variance trade-off between the two estimators.

Like in de Chaisemartin and D'Haultfoeuille (2020), it is straightforward to propose a placebo version of the DID_M^f estimator that one can use to test Assumption 4. To do so, one just needs to replace $Y_{g,t} - Y_{g,t-1}$ by $Y_{g,t-1} - Y_{g,t-2}$ in Equations (14) and (15) above, and exclude from the estimation groups experiencing a change in any of their treatments from $t - 2$ to $t - 1$. The resulting placebo estimator compares the outcome evolution of switchers and non-switchers, before switchers switch.

The DID_M^f estimator can be extended to accommodate discrete non-binary treatments taking values in $\mathcal{D}_1 = \{0, \dots, \bar{d}\}$, like the DID_M estimator in de Chaisemartin and D'Haultfoeuille (2020) (see Section 4 of the Web Appendix of de Chaisemartin and D'Haultfoeuille, 2020). For all

$t \in \{2, \dots, T\}$, for all $(d, d') \in (\mathcal{D}_1)^2$, and for all $d^{-1} \in \mathcal{D}_2 \times \dots \times \mathcal{D}_K$, let

$$\text{DID}_{d,d',d^{-1},t}^f = [1\{d' < d\} - 1\{d < d'\}] \left[\sum_{g \in \mathcal{G}_{d,d',d^{-1},t}} \frac{N_{g,t}}{N_{d,d',d^{-1},t}} [Y_{g,t} - Y_{g,t-1}] \right. \\ \left. - \sum_{g \in \mathcal{G}_{d',d',d^{-1},t}} \frac{N_{g,t}}{N_{d',d',d^{-1},t}} [Y_{g,t} - Y_{g,t-1}] \right]$$

be a DID estimator comparing the $t - 1$ -to- t outcome evolution in groups whose first treatment changes from d' to d and whose other treatments are equal to d^{-1} at both dates, to the same outcome evolution in groups whose treatments do not change and with the same treatments in $t - 1$. With a non-binary treatment, the DID_M^f estimator is a weighted average of the $\text{DID}_{d,d',d^{-1},t}^f$ estimators, across d, d', d^{-1} , and t , normalized by the average change of the first treatment among switchers, to ensure the estimator can be interpreted as an effect produced by a one-unit increase of the first treatment.

Similarly, under Assumption 5, and getting back to the binary treatment case, δ_2 can be unbiasedly estimated by a weighted average of DID's. For all $t \in \{1, \dots, T - 1\}$, for all $(d, d') \in (\mathcal{D}_1)^2$, and for all $d^{-1} \in \mathcal{D}_2 \times \dots \times \mathcal{D}_K$, let $N_{d,d',d^{-1},t+1,t} = \sum_{g \in \mathcal{G}_{d,d',d^{-1},t+1,t}} N_{g,t}$ denote the total population, at period t , of groups in $\mathcal{G}_{d,d',d^{-1},t+1}$. Then, let

$$\text{DID}_{+,d^{-1},t}^b = \sum_{g \in \mathcal{G}_{0,1,d^{-1},t+1,t}} \frac{N_{g,t}}{N_{0,1,d^{-1},t+1,t}} (Y_{g,t} - Y_{g,t+1}) - \sum_{g \in \mathcal{G}_{0,0,d^{-1},t+1,t}} \frac{N_{g,t}}{N_{0,0,d^{-1},t+1,t}} (Y_{g,t} - Y_{g,t+1}), \\ \text{DID}_{-,d^{-1},t}^b = \sum_{g \in \mathcal{G}_{1,1,d^{-1},t+1,t}} \frac{N_{g,t}}{N_{1,1,d^{-1},t+1,t}} (Y_{g,t} - Y_{g,t+1}) - \sum_{g \in \mathcal{G}_{1,0,d^{-1},t+1,t}} \frac{N_{g,t}}{N_{1,0,d^{-1},t+1,t}} (Y_{g,t} - Y_{g,t+1}).$$

In contrast to $\text{DID}_{+,d^{-1},t}^f$, which is a ‘‘forward’’ DID, $\text{DID}_{+,d^{-1},t}^b$ is a ‘‘backward’’ DID, from the future to the past. It compares the $t + 1$ -to- t outcome evolution of groups whose first treatment goes from 0 to 1 from $t + 1$ to t while their other treatments are equal to d^{-1} at both dates, to the outcome evolution of groups whose first and other treatments are respectively equal to 0 and d^{-1} at both dates. $\text{DID}_{-,d^{-1},t}^b$ has a similar interpretation, except that it compares groups whose first treatment is equal to 1 at both dates to groups whose first treatment goes from 1 to 0 from $t + 1$ to t . Let

$$\text{DID}_M^b = \sum_{t=1}^{T-1} \sum_{d^{-1} \in \{0,1\}^{K-1}} \left(\frac{N_{0,1,d^{-1},t+1,t}}{N_{\mathcal{S}_2}} \text{DID}_{+,d^{-1},t}^b + \frac{N_{1,0,d^{-1},t+1,t}}{N_{\mathcal{S}_2}} \text{DID}_{-,d^{-1},t}^b \right) \quad (17)$$

if $N_{\mathcal{S}_2} > 0$, and $\text{DID}_M^b = 0$ if $N_{\mathcal{S}_2} = 0$. One can show that if Assumptions 1-2 and 5 hold, $E[\text{DID}_M^b] = \delta_2$.

4.4 Dynamic treatment effects

With a single treatment $D_{g,t}^s$, DID_M^f can be used to estimate the effect of the current value of $D_{g,t}^s$, allowing for dynamic effects. Assume that $(D_{g,t}^1, \dots, D_{g,t}^K) = (D_{g,t}^s, \dots, D_{g,t-(K-1)}^s)$. Then,

our potential outcome notation allows the current treatment and its first $K - 1$ lags to affect the outcome, so DID_M^f is an estimator of the effect of the current value of $D_{g,t}^s$ robust to dynamic effects up to $K - 1$ lags. This is an improvement over the DID_M estimator in de Chaisemartin and D’Haultfoeuille (2020), which is not robust to dynamic effects, except with a binary and staggered treatment. To achieve some robustness to dynamic effects, DID_M^f restricts the estimation to groups that did not experience a treatment change from $t - K$ to $t - 1$. For instance, with $K = 2$ and $(D_{g,t}^1, D_{g,t}^2) = (D_{g,t}^s, D_{g,t-1}^s)$, the DID_M^f estimator compares groups with $(D_{g,t-2}^s, D_{g,t-1}^s, D_{g,t}^s) = (0, 0, 1)$ to groups with $(D_{g,t-2}^s, D_{g,t-1}^s, D_{g,t}^s) = (0, 0, 0)$, and groups with $(D_{g,t-2}^s, D_{g,t-1}^s, D_{g,t}^s) = (1, 1, 0)$ to groups with $(D_{g,t-2}^s, D_{g,t-1}^s, D_{g,t}^s) = (1, 1, 1)$. On the other hand, the DID_M^f estimator may not be used to estimate the effect of past treatments on the outcome. For instance, with $K = 2$ and $(D_{g,t}^1, D_{g,t}^2) = (D_{g,t-1}^s, D_{g,t}^s)$, \mathcal{S}_1 is empty: for any group g such that $(D_{g,t-2}^s \neq D_{g,t-1}^s = D_{g,t}^s)$, there cannot exist another group g' such that $(D_{g',t-2}^s = D_{g',t-1}^s = D_{g',t}^s)$, and $(D_{g',t-2}^s, D_{g',t-1}^s) = (D_{g,t-2}^s, D_{g,t-1}^s)$.

The opposite applies to DID_M^b : it may not be used to estimate the effect of the current treatment allowing for dynamic effects, but it may be used to estimate the effect of past treatments on the outcome. For instance, with $K = 2$ and $(D_{g,t}^1, D_{g,t}^2) = (D_{g,t-1}^s, D_{g,t}^s)$, \mathcal{S}_2 is not empty: it contains all (g, t) cells such that $D_{g,t-1}^s \neq D_{g,t}^s = D_{g,t+1}^s$, for which there exists another group g' such that $D_{g',t-1}^s = D_{g',t}^s = D_{g',t+1}^s = D_{g,t+1}^s$. Then, DID_M^b is a weighted average of two types of DID. DIDs of the first type compare the $t + 1$ to t outcome evolution between groups such that $D_{g,t-1}^s = 1, D_{g,t}^s = 0, D_{g,t+1}^s = 0$ and groups such that $D_{g,t-1}^s = 0, D_{g,t}^s = 0, D_{g,t+1}^s = 0$. DIDs of the second type compare the $t + 1$ to t outcome evolution between groups such that $D_{g,t-1}^s = 1, D_{g,t}^s = 1, D_{g,t+1}^s = 1$ and groups such that $D_{g,t-1}^s = 0, D_{g,t}^s = 1, D_{g,t+1}^s = 1$. If the current outcome only depends on the current treatment and its first lag, and if Assumption 5 holds for $(D_{g,t}^1, D_{g,t}^2) = (D_{g,t-1}^s, D_{g,t}^s)$, then DID_M^b is unbiased for the average effect of switching the treatment’s first lag from 0 to 1 holding the current treatment fixed, across all (g, t) s in \mathcal{S}_2 .

Of course, assuming that the current outcome only depends on the current treatment and its first lag is restrictive. One could instead assume, say, that the current outcome only depends on the current treatment and its first two lags. Then, with $K = 3$ and $(D_{g,t}^1, D_{g,t}^2, D_{g,t}^3) = (D_{g,t-2}^s, D_{g,t-1}^s, D_{g,t}^s)$, DID_M^b is unbiased for the average effect of switching the treatment’s second lag from 0 to 1, holding the current treatment and its first lag fixed, across all (g, t) cells in \mathcal{S}_2 . \mathcal{S}_2 now becomes the set of all (g, t) cells such that $D_{g,t-2}^s \neq D_{g,t-1}^s = D_{g,t}^s = D_{g,t+1}^s$ and for which there exists another group g' such that $D_{g',t-2}^s = D_{g',t-1}^s = D_{g',t}^s = D_{g',t+1}^s = D_{g,t+1}^s$. \mathcal{S}_2 contains fewer cells with $K = 3$ and $(D_{g,t}^1, D_{g,t}^2, D_{g,t}^3) = (D_{g,t-2}^s, D_{g,t-1}^s, D_{g,t}^s)$ than with $K = 2$ and $(D_{g,t}^1, D_{g,t}^2) = (D_{g,t-1}^s, D_{g,t}^s)$: allowing more treatment lags to affect the outcome may be more plausible, but it may also result in less precise estimators, that apply to a smaller population. Note also that with dynamic effects up to $K - 1$ treatment lags, DID_M^b can be used to estimate the effect of the $K - 1$ th lag, but it cannot be used to estimate the effect of earlier lags.

Overall, DID_M^f and DID_M^b can be used for some but not for all purposes in the presence of a single treatment with dynamic effects. Assuming constant treatment effects, one can use a TWFE regression of the outcome on the treatment and its lags, the so-called distributed lag regression, to separately estimate the effect of the current and past treatments on the outcome. Separately estimating each of those effects while allowing for heterogeneous treatment effects is inherently difficult. This may be the reason why a substantial branch of the heterogeneity-robust DID literature has instead proposed to estimate the total effect of current and past treatments on the outcome (see Callaway and Sant’Anna, 2021; Sun and Abraham, 2021; de Chaisemartin and D’Haultfœuille, 2021a; Borusyak, Jaravel and Spiess, 2021). This literature has focused on the case with one treatment. In Section 2 of the Web Appendix, we extend that literature, and propose estimators of instantaneous and dynamic effects when there are several binary and staggered treatments. Those estimators can also be used in the presence of a single treatment that can change multiple times, to isolate the effect of each treatment change. For instance, with a treatment that can switch on and then off, one may be interested in separately estimating the effect of switching the treatment on/off.

4.5 Inference

In Section 3 of the Web Appendix, we prove the asymptotic normality of DID_M^f when the number of groups goes to infinity, and we propose confidence intervals. Those results are established under similar assumptions and arguments as those used to show the asymptotic normality of the DID_M estimator in de Chaisemartin and D’Haultfœuille (2020) (see Theorem S6 in the Web Appendix therein), without any important conceptual difference. One limitation of this approach, though, is that the asymptotic approximation may not be accurate. DID_M^f compares carefully selected treatment and control groups, and it could be the case that only a small number of groups can be included in those comparisons. The larger the number of treatments, the more likely it is that DID_M^f uses data from a small number of groups. In this section, we deal with this issue by proposing confidence intervals that are exact in a finite sample of groups under a normality assumption, in the spirit of Donald and Lang (2007). The exactness of those confidence intervals relies on strong conditions, but they remain asymptotically valid under much weaker assumptions. The main price to pay for using them, rather than those described in Section 3 of the Web Appendix, is that doing so may result in an adjustment of the definition of δ_1 , as explained below.

To ease the exposition, in this section we condition on \mathbf{D} . Accordingly, functions of \mathbf{D} can be treated as non-stochastic terms. For simplicity, we also assume that $N_{g,t} = 1$ for all (g, t) . Let $s = 1, \dots, S$ index “switches”, that is to say a $K + 2$ -uple (d, d', d^{-1}, t) with $d \neq d'$ for which there exists (g, g') satisfying $D_{g,t}^1 = d$, $D_{g,t-1}^1 = D_{g',t-1}^1 = D_{g',t}^1 = d'$ and $D_{g,t}^{-1} = D_{g,t-1}^{-1} = D_{g',t}^{-1} =$

$D_{g',t-1}^{-1} = d^{-1}$. Then, note that

$$\text{DID}_M = \sum_{s=1}^S \alpha_s \text{DID}_s,$$

for some non-stochastic weights $(\alpha_s)_{s=1,\dots,S}$, where DID_s is the DID corresponding to switch s . Let \mathcal{G}_s denote the set of groups intervening in DID_s , either as a “switcher” or as a “control”. The following assumption is needed to ensure the validity of our approach.

Assumption 6 (Non-overlapping groups) *for any $(s, s') \in \{1, \dots, S\}^2$, $s \neq s'$, $\mathcal{G}_s \cap \mathcal{G}_{s'} = \emptyset$.*

Note that Assumption 6 automatically holds with $T = 2$. Otherwise, it is more likely to hold if T is small. When Assumption 6 fails, we can ensure it holds on a modified sample, by removing groups belonging to several sets \mathcal{G}_s from all those sets except one. This sample modification will modify the estimator DID_M^f . It may also lead to removing a switching group from a set \mathcal{G}_s . This would change the target parameter, which would become the average treatment effect across all switching cells in the modified sample, in lieu of δ_1 , the average treatment effect across all switching cells in the original sample. For simplicity, we still denote the parameter and its estimator on the modified sample δ_1 and DID_M^f .

Our confidence interval relies on the following variance estimator:

$$\widehat{V} = \sum_{s=1}^S \alpha_s^2 \left[\frac{1}{n_{1s}(n_{1s} - 1)} \sum_{g \in \mathcal{G}_{1s}} (\Delta Y_{g,t_s} - \overline{\Delta Y}_{1s})^2 + \frac{1}{n_{0s}(n_{0s} - 1)} \sum_{g \in \mathcal{G}_{0s}} (\Delta Y_{g,t_s} - \overline{\Delta Y}_{0s})^2 \right],$$

where \mathcal{G}_{1s} (resp. \mathcal{G}_{0s}) is the subset of switching (resp. control) cells in \mathcal{G}_s , $n_{ks} = \text{card}(\mathcal{G}_{ks})$, and $\overline{\Delta Y}_{ks}$ is the average of $\Delta Y_{g,t_s}$ over $g \in \mathcal{G}_{ks}$. Our definition of \widehat{V} uses the convention that $0/0=0$.

Next, let $q_{1-\alpha}$ denote the quantile of order $1 - \alpha$ of $|T|$, defined as

$$T = \left(\frac{\sum_{s=1}^S \alpha_s^2 (1/n_{1s} + 1/n_{0s})}{\sum_{s=1}^S \alpha_s^2 [W_{1s}/[n_{1s}(n_{1s} - 1)] + W_{0s}/[n_{0s}(n_{0s} - 1)]]} \right)^{1/2} \times Z, \quad (18)$$

where $(Z, W_{01}, W_{11}, \dots, W_{0S}, W_{1S})$ are independent of the data, mutually independent and satisfy $Z \sim \mathcal{N}(0, 1)$ and $W_{ks} \sim \chi^2(n_{ks} - 1)$. Note that $q_{1-\alpha}$ does not have a closed-form expression, but it can be approximated by simulations.

Then, the confidence interval of order $1 - \alpha$ we consider is

$$\text{CI}_{1-\alpha}^{\text{ex}} = \left[\text{DID}_M \pm q_{1-\alpha} \sqrt{\widehat{V}} \right].$$

Below, we introduce two assumptions under which $\text{CI}_{1-\alpha}^{\text{ex}}$ is valid: under Assumption 7, $\text{CI}_{1-\alpha}^{\text{ex}}$ is valid in finite samples; under Assumption 8, $\text{CI}_{1-\alpha}^{\text{ex}}$ is valid asymptotically.

Assumption 7 (Restrictions for finite-sample validity of $\text{CI}_{1-\alpha}^{\text{ex}}$)

1. For all $s = 1, \dots, S$ and $g \in \mathcal{G}_s$, $Y_{g,t_s}(d_s^1, d_s^{-1}) - Y_{g,t_s}(d_s^{1'}, d_s^{-1}) = \delta_{1s}$ where $(t_s, d_s^1, d_s^{1'}, d_s^{-1})$ are the $(t, d^1, d^{1'}, d^{-1})$ associated with s and δ_{1s} is non-stochastic.
2. For all s and $g \in \mathcal{G}_s$, $\Delta Y_{g,t_s}(0, d_s^{-1}) \sim \mathcal{N}(\mu_s, \sigma^2)$.

Point 1 of Assumption 7 assumes that the first treatment's effect is homogeneous within each set of groups s , but may vary across s . Point 2 of Assumption 7 assumes that $\Delta Y_{g,t_s}(0, d_s^{-1})$ is normally distributed and homoskedastic: the variance of $\Delta Y_{g,t_s}(0, d_s^{-1})$ should not depend on s .

Assumption 8 (Restrictions for asymptotic validity of $CI_{1-\alpha}^{\text{ex}}$)

1. There exists G_0 such that for all $G \geq G_0$, $\mathcal{S}_G := \{(d, d', d^{-1}, t) : N_{d,d',d^{-1},t} > 0, N_{d',d',d^{-1},t} > 0\}$ does not vary across G and is finite. We denote by \bar{S} its cardinal.¹¹
2. For all $(k, s) \in \{0, 1\} \times \{1, \dots, \bar{S}\}$, the $(\Delta Y_{g,t_s})_{g \in \mathcal{G}_{ks}}$ are i.i.d. and for $g \in \mathcal{G}_{ks}$, $E[\Delta Y_{g,t_s}^2] < \infty$ and $V(\Delta Y_{g,t_s}) > 0$.
3. For all $(k, s) \in \{0, 1\} \times \{1, \dots, \bar{S}\}$, $\liminf_{G \rightarrow \infty} n_{ks}/G > 0$.

Assumption 8 does not make any treatment-effect homogeneity or homoscedasticity assumption. Note that we impose that the $(\Delta Y_{g,t_s})_{g \in \mathcal{G}_{ks}}$ are identically distributed for simplicity. If we instead assumed that the $(\Delta Y_{g,t_s})_{g \in \mathcal{G}_{ks}}$ are independent but not identically distributed, one could still show that $CI_{1-\alpha}^{\text{ex}}$ is asymptotically conservative under appropriate regularity conditions, as in, e.g., Theorem S6 in the Web Appendix of de Chaisemartin and D'Haultfœuille (2020).

The following theorem shows that $CI_{1-\alpha}^{\text{ex}}$ is exact under Assumption 7, and asymptotically valid under Assumption 8.

Theorem 5 *If Assumptions 1-2, 4, and 6 hold, then, for any $\alpha \in (0, 1)$:*

1. *if Assumption 7 further holds, $P(\delta_1 \in CI_{1-\alpha}^{\text{ex}}) = 1 - \alpha$.*
2. *if Assumption 8 further holds, $\lim_{G \rightarrow \infty} P(\delta_1 \in CI_{1-\alpha}^{\text{ex}}) = 1 - \alpha$.*

Our approach in this section generalizes that in Donald and Lang (2007) to designs with more than two time periods and/or several treatments. A difference with that paper is that our confidence intervals use critical values from a non-standard distribution, instead of critical values from a t-distribution.

¹¹Accordingly, we keep the same indexation for switches $s \in \{1, \dots, \bar{S}\}$ for all $G \geq G_0$.

5 Application

In this section, we revisit Hotz and Xiao (2011).¹² Unfortunately, many tables in this paper rely on proprietary data. The only table with TWFE regressions with several treatments that we can replicate is Table 11. Therefore, we focus on this table in our replication, though it is not the paper’s main table.

Hotz and Xiao (2011) use a panel of the 50 US states and the District of Columbia, in 1987, 1992, and 1997, to estimate the effect of state center-based daycare regulations, namely the minimum years of schooling required to be the director of a center-based care and the minimum staff-to-child ratio, on the demand for family home daycare. Family home day cares are not subject to those regulations. More stringent regulations may increase the cost of center-based establishments, but may also increase their safety and quality. Accordingly, the effects of those regulations on the demand for family home daycare is ambiguous. The distributions of these regulations are shown in Table 1. The minimum years of schooling is a discrete treatment taking six values included between 0 (no minimum) and 16, with 14 (associate degree) being the most frequent value. The minimum staff-to-child ratio is a also discrete treatment variable, taking seven values included between 0 (no minimum) and 1/3 (one professional per three children), with 1/4 being the most frequent value.

¹²This paper is the only one, in the census of TWFE papers published by the AER from 2010 to 2012 that we conducted in de Chaisemartin and D’Haultfœuille (2020), that has several treatments in the regression, relies at least partially on non-proprietary data, and for which the treatments are not continuous (thus making it possible to compute the DID_M^f estimator).

Table 1: Distribution of the two treatments in Hotz and Xiao (2011)

Min. years of schooling	# of (g,t) cells
0	26
12	36
12.5	5
13	4
14	61
16	21
Min. staff-to-child ratio	# of (g,t) cells
0	5
1/8	2
1/7	4
1/6	30
1/5	21
1/4	82
1/3	9

Hotz and Xiao (2011) regress the revenue of family home day cares in state g and year t on state fixed effects, year fixed effects, 12 control variables, the minimum years of schooling required to be the director of a center-based care, the minimum staff-to-child ratio, and two indicators for whether there is no such minima, to allow for potentially non-linear effects. In Column (3) of their Table 11, the coefficient on the minimum years of schooling treatment, $\hat{\beta}_{fe}^X$, is equal to -0.445 and is highly significant (95% confidence interval= $[-0.735, -0.155]$),¹³ thus suggesting that increasing by one the years of schooling required for directors of center-based daycare decreases the revenue of family home daycare by 0.44 million USD.

Dropping the 12 control variables from the regression does not affect that conclusion very much: the coefficient on the minimum years of schooling treatment, $\hat{\beta}_{fe}$, is now equal to -0.566 and is still highly significant (95% confidence interval= $[-0.852, -0.280]$). Below, we study $\hat{\beta}_{fe}$, rather than $\hat{\beta}_{fe}^X$, the coefficient estimated by Hotz and Xiao (2011). This is to ensure that the TWFE estimator we study is comparable to the DID_M^f estimator we compute below: while the DID_M^f estimator can be extended to allow for control variables, the sample on which it is computed in this application is not large enough to include 12 control variables.

¹³This confidence interval is slightly larger than that in Hotz and Xiao (2011), because we cluster standard errors at the state rather than at the state \times year level, which is more in line with the standard practice in empirical work (see Bertrand, Duflo and Mullainathan, 2004).

We now show that $\widehat{\beta}_{fe}$ may not be robust to heterogeneous effects across state and years, and may also be contaminated by the effects of the other treatments in the regression. Following Corollary 1, this coefficient can be decomposed into the sum of four terms. The first term is a weighted sum of the effects of increasing by one the years of schooling required in 127 state \times year cells, where 44 effects receive a positive weight and 83 receive a negative weight, and where the positive and negative weights respectively sum to 7.897 and -6.897. The second term is a sum of the effects of not having a requirement on directors' years of schooling in 26 state \times year cells, where 11 effects receive a positive weight and 15 receive a negative weight, and where the positive and negative weights respectively sum to 0.148 and -0.148. The third term is a sum of the effects of increasing by one the staff to child ratio in 148 state \times year cells, where 51 effects receive a positive weight and 97 receive a negative weight, and where the positive and negative weights respectively sum to 0.160 and -0.160. The last term is a sum of the effects of not having a requirement on staff to child ratio in 5 state \times year cells, where 4 effects receive a positive weight and 1 receive a negative weight, and where the positive and negative weights respectively sum to 0.055 and -0.055. Results are similar for the other three treatment coefficients in the regression, except that the contamination weights attached to them are even larger. For instance, for the coefficient on the staff to child ratio treatment, the weighted sum of the effects of the minimum years of schooling treatment has positive and negative weights summing to 246.222 and -246.222.

When the other three treatment variables are dropped from the regression, the coefficient on the minimum years of schooling becomes small (-0.020) and insignificant (95% confidence interval= $[-0.114, 0.074]$). We follow Theorem 3 to decompose the coefficient in this "short" regression, and compare it to the coefficient in the "long" regression with the four treatments. The short regression's coefficient can be decomposed into the sum of four terms. The first term is a weighted sum of the effects of increasing by one the years of schooling required in 127 state \times year cells, where 56 cells receive a positive weight and 71 receive a negative weight, and where the positive and negative weights respectively sum to 1.759 and -0.759. Thus, the short regression has considerably smaller negative weights in this first term than the long regression. The second term is a sum of the effects of not having a requirement on directors' years of schooling in 26 state \times year cells, where 5 effects receive a positive weight and 21 receive a negative weight, and where the positive and negative weights respectively sum to 0.008 and -0.077. The third term is a sum of the effects of increasing by one the staff to child ratio in 148 state \times year cells, where 61 effects receive a positive weight and 87 receive a negative weight, and where the positive and negative weights respectively sum to 0.030 and -0.022. The last term is a sum of the effects of not having a requirement on staff to child ratio in 5 state \times year cells, where all effects receive a negative weight, and where the negative weights sum to -0.035. Thus, the short regression also has considerably less contamination weights than the long regression. Accordingly, the estimated maximal bias in Corollary 2 is almost five times lower for the short than for

the long regression ($4.233 \times B$ versus $20.741 \times B$), so the short regression is preferable per this maximal-bias metric.

Finally, we compute the estimator proposed in Section 4, for the minimum years of schooling treatment, controlling for the staff-to-child ratio treatment. Our estimators do not assume linear treatment effects, so we do not need to control for the indicators for whether there is no such minima.

There are 127 (g, t) cells with a non-zero minimum years of schooling. On the other hand, there are only five (g, t) cells in \mathcal{S}_1 , all of which have a non-zero minimum years of schooling. The 5 (g, t) cells our estimator applies to are (Kentucky,1992), (Minnesota,1992), (Utah,1992), (Vermont,1992), and (Rhode Island,1997).¹⁴ Of the 122 (g, t) cells we lose when focusing on \mathcal{S}_1 , 93 belong to the first subgroup in Appendix 1, 19 belong to the second or third subgroup, and 10 belong to the fourth or fifth subgroup. Therefore, the vast majority of the cells we lose do not experience any change of their minimum years of schooling, so their treatment effect cannot be identified under a parallel trends assumption. We may seem to lose 19 cells by imposing only a minimal parallel trends assumption. In reality, estimating the treatment effects of 14 of the 19 cells in the second or third subgroup would also require assuming that the effect of the minimum staff-to-child ratio is homogeneous: either their minimum staff-to-child ratio also changes when their minimum years of schooling changes, or they cannot be matched to a control state with the same baseline treatments.

We find that $\text{DID}_M^f = -0.029$. DID_M^f uses data from 5 switching and 19 control (g, t) cells, so the asymptotic approximation in Section 3 of the Web Appendix may not be very reliable for that estimator. Instead, we compute the confidence interval $\text{CI}_{1-\alpha}^{\text{ex}}$ for $\alpha = 0.95$ and find that it is equal to $[-0.821, 0.807]$.¹⁵ In this application, the assumption that the first-differenced outcome is normally distributed is not rejected. We conduct a Shapiro-Wilk test separately for the 1987 to 1992 and for the 1992 to 1997 first differences, as the test assumes independent observations. None of the two tests is rejected (p-value= 0.98 and 0.46, respectively).

To gain precision, one may further impose Assumption 5. Doing so allows us to use DID_M^b to estimate the treatment effect in five (g, t) cells in \mathcal{S}_2 . \mathcal{S}_1 and \mathcal{S}_2 do not overlap and have the same numbers of cells, so we can also use $1/2(\text{DID}_M^f + \text{DID}_M^b)$ to estimate δ , the average treatment effect in $\mathcal{S}_1 \cup \mathcal{S}_2$. We find that $1/2(\text{DID}_M^f + \text{DID}_M^b) = -0.016$. $1/2(\text{DID}_M^f + \text{DID}_M^b)$ uses data from 50 (g, t) cells, coming from 30 different states. The asymptotic approximation

¹⁴For the staff-to-child ratio treatment, the set \mathcal{S}_1 is even smaller as it only contains two (g, t) cells. This is why we focus on the minimum-years-of-schooling treatment.

¹⁵ $\text{CI}_{1-\alpha}^{\text{ex}}$ relies on Assumption 6, which does not hold in our data: Rhode Island and Washington are used twice in DID_M^f . Removing these two states in one of the two s they belong to (using the notation in Section 4.5) changes very slightly the value of DID_M^f (-0.0072 in lieu of -0.029).

in Section 3 of the Web Appendix may be more reasonable for that estimator,¹⁶ so we follow Theorem 7 therein to compute a 95% confidence interval for δ . We find that this confidence interval is equal to $[-0.126, 0.094]$. We also test the equality between δ and β_{fe} , and reject the null hypothesis at all conventional levels (p-value= 4×10^{-4}). Hence, as discussed above, we can reject the hypothesis that the effects of the minimum years of schooling and staff-to-child ratio treatments are homogenous.

Let us summarize our results. Using a TWFE regression with several treatments, Hotz and Xiao (2011) find that increasing the years of schooling required for directors of center-based daycare significantly decreases the revenue of family home daycare. We show that in the presence of heterogeneous treatment effects, their regression estimates a highly-non-convex combination of the effects of the years of schooling treatment, and is contaminated by the effects of the other treatments. Therefore, their finding may not be robust to heterogeneous treatment effects. Then, we use our robust estimators to assess if, in the presence of heterogeneous effects, one can conclude, for at least a subset of (g, t) cells, that increasing the years of schooling requirement significantly decreases the revenue of family home daycare. The answer is negative, as our estimators are insignificant. Moreover, one of our estimators is significantly different from the TWFE estimator, thus allowing us to reject the null hypothesis that the effects of all treatments are constant in this application. Overall, there is no evidence that the finding in Hotz and Xiao (2011) is robust to heterogeneous effects, while there is evidence that treatment effects are heterogeneous in this application.

Table 2: Estimators of the effect of the minimum years of schooling treatment

	Estimate	95% Confidence Interval
$\widehat{\beta}_{fe}^X$	-0.445	$[-0.735, -0.155]$
$\widehat{\beta}_{fe}$	-0.566	$[-0.852, -0.280]$
$\widehat{\beta}_s$	-0.022	$[-0.114, 0.074]$
DID_M^f	-0.029	$[-0.821, 0.807]$
$1/2(\text{DID}_M^f + \text{DID}_M^b)$	-0.016	$[-0.126, 0.094]$

¹⁶To verify that, we considered simulations with the same design as in the application but with no effects of the treatments, and $(\Delta Y_{g,2}(\mathbf{0}), \Delta Y_{g,3}(\mathbf{0}))$ drawn either from a normal distribution $\mathcal{N}(\mathbf{0}, \Sigma)$, with Σ equal to the estimated variance matrix on the sample, or from the empirical distribution of $(\Delta Y_{g,2}, \Delta Y_{g,3})$. In both cases, the coverage of our confidence interval was higher than 95% (95.4% and 99.3%, respectively).

6 Conclusion

In this paper, we show that treatment coefficients in TWFE regressions with several treatments may not be robust to heterogeneous effects, and could be contaminated by the effects of other treatments in the regression. We propose alternative DID estimators that are robust to heterogeneous effects and do not suffer from this contamination problem.

In most instances where TWFE and DID estimators are used, it is likely that besides the main treatment of interest, many other determinants of the outcome change over the study period. We show that in the presence of heterogeneous treatment effects, failing to control for those other treatments, be it in a TWFE regression or using an heterogeneity-robust DID estimator, may lead to a biased estimator, even if those other treatments are uncorrelated with the main treatment of interest. Accordingly, all those other treatments should be controlled for, but our results also show that a non-parametric DID estimator robust to the heterogeneous effects of many treatments will often be subject to a curse of dimensionality. Data-driven methods to select the treatments that should be controlled for, as well as more parametric methods to control for them, would be useful additions to the econometrics literature. In the meantime, applied researchers could discuss more systematically whether other treatments than the one under consideration have changed over their study period. If so, they could assess if their estimates are robust to controlling for at least some of those other treatments, using the tools provided in this paper.

References

- Abadie, Alberto.** 2005. “Semiparametric Difference-in-Differences Estimators.” *Review of Economic Studies*, 72(1): 1–19.
- Ashenfelter, Orley.** 1978. “Estimating the effect of training programs on earnings.” *The Review of Economics and Statistics*, 47–57.
- Bertrand, Marianne, Esther Dufo, and Sendhil Mullainathan.** 2004. “How much should we trust differences-in-differences estimates?” *The Quarterly Journal of Economics*, 119(1): 249–275.
- Bojinov, Iavor, Ashesh Rambachan, and Neil Shephard.** 2021. “Panel experiments and dynamic causal effects: A finite population perspective.” *Quantitative Economics*, 12(4): 1171–1196.
- Borusyak, Kirill, and Xavier Jaravel.** 2017. “Revisiting event study designs.” Working Paper.
- Borusyak, Kirill, Xavier Jaravel, and Jann Spiess.** 2021. “Revisiting event study designs: Robust and efficient estimation.” *arXiv preprint arXiv:2108.12419*.
- Callaway, Brantly, and Pedro H.C. Sant’Anna.** 2021. “Difference-in-Differences with Multiple Time Periods.” *Journal of Econometrics*, 225: 200–230.
- de Chaisemartin, C, and X D’Haultfœuille.** 2018. “Fuzzy Differences-in-Differences.” *The Review of Economic Studies*, 85(2): 999–1028.
- de Chaisemartin, Clement, and Xavier D’Haultfœuille.** 2020. “Two-way fixed effects estimators with heterogeneous treatment effects.” *American Economic Review*, 110(9): 2964–96.
- de Chaisemartin, Clément, and Xavier D’Haultfœuille.** 2021*a*. “Difference-in-Differences Estimators of Intertemporal Treatment Effects.” arXiv preprint arXiv:2007.04267.
- de Chaisemartin, Clément, and Xavier d’Haultfoeuille.** 2021*b*. “Two-way fixed effects regressions with several treatments.” *arXiv preprint arXiv:2012.10077, v4*.
- D’Haultfœuille, Xavier, and Purevdorj Tuvaandorj.** 2022. “A Robust Permutation Test for Subvector Inference in Linear Regressions.” arXiv preprint 2205.06713.
- Donald, Stephen G, and Kevin Lang.** 2007. “Inference with difference-in-differences and other panel data.” *The review of Economics and Statistics*, 89(2): 221–233.

- Goldsmith-Pinkham, Paul, Peter Hull, and Michal Kolesár.** 2021. “On Estimating Multiple Treatment Effects with Regression.” arXiv preprint arXiv:2106.05024.
- Goodman-Bacon, Andrew.** 2021. “Difference-in-differences with variation in treatment timing.” *Journal of Econometrics*, 225: 254–277.
- Holland, Paul W.** 1986. “Statistics and causal inference.” *Journal of the American statistical Association*, 81(396): 945–960.
- Holland, Paul W, and Donald B Rubin.** 1987. “Causal inference in retrospective studies.” *ETS Research Report Series*, 1987(1): 203–231.
- Hotz, V Joseph, and Mo Xiao.** 2011. “The impact of regulations on the supply and quality of care in child care markets.” *American Economic Review*, 101(5): 1775–1805.
- Hull, Peter.** 2018. “Estimating Treatment Effects in Mover Designs.” arXiv preprint 1804.06721.
- Meinhofer, Angélica, Allison Witman, Jesse Hinde, and Kosali Simon.** 2021. “Marijuana liberalization policies and perinatal health.” *Journal of Health Economics*, 102537.
- Robins, James.** 1986. “A new approach to causal inference in mortality studies with a sustained exposure period-application to control of the healthy worker survivor effect.” *Mathematical modelling*, 7(9-12): 1393–1512.
- Sun, Liyang, and Sarah Abraham.** 2021. “Estimating dynamic treatment effects in event studies with heterogeneous treatment effects.” *Journal of Econometrics*, 225: 175–199.

A Proofs

A.1 Theorem 1

The result directly follows from Theorem 2. If $K = 2$, $D_{g,t}^{-1} = D_{g,t}^2$. Then, $D_{g,t}^{-1} \neq \mathbf{0}^{-1}$ if and only if $D_{g,t}^2 = 1$, and one then has $D_{g,t}^2 \Delta_{g,t}^{-1} = D_{g,t}^2 \Delta_{g,t}^2$.

A.2 Theorem 2

We first establish the following lemma.

Lemma 1 *If Assumptions 1-3 hold, for all $(g, g', t, t') \in \{1, \dots, G\}^2 \times \{1, \dots, T\}^2$,*

$$\begin{aligned} & E(Y_{g,t} | \mathbf{D}) - E(Y_{g,t'} | \mathbf{D}) - (E(Y_{g',t} | \mathbf{D}) - E(Y_{g',t'} | \mathbf{D})) \\ &= D_{g,t}^1 E(\Delta_{g,t}^1 | \mathbf{D}) + E(\Delta_{g,t}^{-1} | \mathbf{D}) - D_{g',t}^1 E(\Delta_{g',t}^1 (D_{g',t}^{-1}) | \mathbf{D}) - E(\Delta_{g',t}^{-1} | \mathbf{D}) \\ & - D_{g,t'}^1 E(\Delta_{g,t'}^1 (D_{g,t'}^{-1}) | \mathbf{D}) - E(\Delta_{g,t'}^{-1} | \mathbf{D}) + D_{g',t'}^1 E(\Delta_{g',t'}^1 (D_{g',t'}^{-1}) | \mathbf{D}) + E(\Delta_{g',t'}^{-1} | \mathbf{D}). \end{aligned}$$

Proof of Lemma 1

For all $(g, t) \in \{1, \dots, G\} \times \{1, \dots, T\}$,

$$\begin{aligned} E(Y_{g,t} | \mathbf{D}) &= E\left(Y_{g,t}(0, \mathbf{0}^{-1}) + D_{g,t}^1(Y_{g,t}(1, D_{g,t}^{-1}) - Y_{g,t}(0, D_{g,t}^{-1}) + Y_{g,t}(0, D_{g,t}^{-1}) - Y_{g,t}(0, \mathbf{0}^{-1}))\right. \\ & \quad \left. + (1 - D_{g,t}^1)(Y_{g,t}(0, D_{g,t}^{-1}) - Y_{g,t}(0, \mathbf{0}^{-1})) \middle| \mathbf{D}\right) \\ &= E\left(Y_{g,t}(0, \mathbf{0}^{-1}) \middle| \mathbf{D}\right) + D_{g,t}^1 E\left(\Delta_{g,t}^1 \middle| \mathbf{D}\right) + E\left(\Delta_{g,t}^{-1} \middle| \mathbf{D}\right) \\ &= E\left(Y_{g,t}(0, \mathbf{0}^{-1}) \middle| \mathbf{D}_g\right) + D_{g,t}^1 E\left(\Delta_{g,t}^1 \middle| \mathbf{D}\right) + E\left(\Delta_{g,t}^{-1} \middle| \mathbf{D}\right), \end{aligned} \tag{19}$$

where the last equality follows from Assumption 2. Moreover, by Assumption 3

$$\begin{aligned} & E\left(Y_{g,t}(0, \mathbf{0}^{-1}) \middle| \mathbf{D}_g\right) - E\left(Y_{g,t'}(0, \mathbf{0}^{-1}) \middle| \mathbf{D}_g\right) - E\left(Y_{g',t}(0, \mathbf{0}^{-1}) \middle| \mathbf{D}_g\right) + E\left(Y_{g',t'}(0, \mathbf{0}^{-1}) \middle| \mathbf{D}_g\right) \\ &= 0. \end{aligned} \tag{20}$$

The result follows by combining (19) and (20).

Proof of Theorem 2

It follows from the Frisch-Waugh theorem and the definition of $\varepsilon_{g,t}$ that

$$E\left(\widehat{\beta}_{fe} \middle| \mathbf{D}\right) = \frac{\sum_{g,t} N_{g,t} \varepsilon_{g,t} E(Y_{g,t} | \mathbf{D})}{\sum_{g,t} N_{g,t} \varepsilon_{g,t} D_{g,t}^1}. \tag{21}$$

Now, by definition of $\varepsilon_{g,t}$ again,

$$\sum_{t=1}^T N_{g,t} \varepsilon_{g,t} = 0 \text{ for all } g \in \{1, \dots, G\}, \quad (22)$$

$$\sum_{g=1}^G N_{g,t} \varepsilon_{g,t} = 0 \text{ for all } t \in \{1, \dots, T\}, . \quad (23)$$

Then,

$$\begin{aligned} & \sum_{g,t} N_{g,t} \varepsilon_{g,t} E(Y_{g,t} | \mathbf{D}) \\ &= \sum_{g,t} N_{g,t} \varepsilon_{g,t} (E(Y_{g,t} | \mathbf{D}) - E(Y_{g,1} | \mathbf{D}) - E(Y_{1,t} | \mathbf{D}) + E(Y_{1,1} | \mathbf{D})) \\ &= \sum_{g,t} N_{g,t} \varepsilon_{g,t} (D_{g,t}^1 E(\Delta_{g,t}^1 | \mathbf{D}) + E(\Delta_{g,t}^{-1} | \mathbf{D}) - D_{1,t}^1 E(\Delta_{1,t}^1 (D_{1,t}^{-1}) | \mathbf{D}) - E(\Delta_{1,t}^{-1} | \mathbf{D})) \\ &\quad - D_{g,1}^1 E(\Delta_{g,1}^1 (D_{g,1}^{-1}) | \mathbf{D}) - E(\Delta_{g,1}^{-1} | \mathbf{D}) + D_{1,1}^1 E(\Delta_{1,1}^1 (D_{1,1}^{-1}) | \mathbf{D}) + E(\Delta_{1,1}^{-1} | \mathbf{D})) \\ &= \sum_{g,t} N_{g,t} \varepsilon_{g,t} (D_{g,t}^1 E(\Delta_{g,t}^1 | \mathbf{D}) + E(\Delta_{g,t}^{-1} | \mathbf{D})) \\ &= \sum_{(g,t): D_{g,t}^1=1} N_{g,t} \varepsilon_{g,t} E(\Delta_{g,t}^1 | \mathbf{D}) + \sum_{(g,t): D_{g,t}^{-1} \neq \mathbf{0}^{-1}} N_{g,t} \varepsilon_{g,t} E(\Delta_{g,t}^{-1} | \mathbf{D}). \end{aligned} \quad (24)$$

The first and third equalities follow from Equations (22) and (23). The second equality follows from Lemma 1. The fourth equality follows from the fact that $\Delta_{g,t}^0(\mathbf{0}^{-1}) = 0$. Finally,

$$\sum_{g,t} N_{g,t} \varepsilon_{g,t} D_{g,t}^1 = \sum_{(g,t): D_{g,t}^1=1} N_{g,t} \varepsilon_{g,t}. \quad (25)$$

Combining (21), (24), (25) yields

$$E(\widehat{\beta}_{fe} | \mathbf{D}) = \sum_{(g,t): D_{g,t}^1=1} \frac{N_{g,t}}{N_1} w_{g,t} E(\Delta_{g,t}^1 | \mathbf{D}) + \sum_{(g,t): D_{g,t}^{-1} \neq \mathbf{0}^{-1}} \frac{N_{g,t}}{N_1} w_{g,t} E(\Delta_{g,t}^{-1} | \mathbf{D}). \quad (26)$$

Then, the first result follows from the law of iterated expectations. Finally, if $K = 2$ or the treatments are mutually exclusive,

$$\sum_{(g,t): D_{g,t}^{-1} \neq \mathbf{0}^{-1}} N_{g,t} \varepsilon_{g,t} E(\Delta_{g,t}^{-1} | \mathbf{D}) = \sum_{k=2}^K \sum_{(g,t): D_{g,t}^k=1} N_{g,t} \varepsilon_{g,t} E(\Delta_{g,t}^{-1} | \mathbf{D}).$$

Moreover, by definition of $\varepsilon_{g,t}$, $\sum_{(g,t): D_{g,t}^k=1} N_{g,t} \varepsilon_{g,t} = 0$ for all $k = 2, \dots, K-1$. The second result follows.

Proof of Theorem 3

The proof is the same as that of Theorem 1, with just one difference: we do not have $\sum_{(g,t): D_{g,t}^2=1} N_{g,t} \times \varepsilon_{g,t}^s = 0$, since $\varepsilon_{g,t}^s$ is not orthogonal to $D_{g,t}^2$ in general.

Proof of Corollary 2

The result directly follows from Theorems 1 and 3, the triangle inequality, and the fact there is a real number B such that $|\Delta_{g,t}^1| \leq B$ and $|\Delta_{g,t}^2| \leq B$ for all (g, t) . The first bound is reached when $\Delta_{g,t}^1 = B \times (2 \times 1\{w_{g,t} \geq 1\} - 1)$ and $\Delta_{g,t}^2 = B(2 \times 1\{w_{g,t} \geq 0\} - 1)$, the second bound is reached when $\Delta_{g,t}^1 = B \times (2 \times 1\{w_{g,t}^s \geq 1\} - 1)$ and $\Delta_{g,t}^2 = B(2 \times 1\{w_{g,t}^s \geq 0\} - 1)$.

Theorem 4

First, by definition of DID_M^f ,

$$\text{DID}_M^f = \sum_{t=2}^T \sum_{d^{-1} \in \{0,1\}^{K-1}} \frac{N_{1,0,d^{-1},t}}{N_{\mathcal{S}_1}} \text{DID}_{+,d^{-1},t}^f + \frac{N_{0,1,d^{-1},t}}{N_{\mathcal{S}_1}} \text{DID}_{-,d^{-1},t}^f, \quad (27)$$

using here the convention that $0/0 = 0$. Let $t \geq 2$ and $d^{-1} \in \{0,1\}^{K-1}$ be such that $N_{1,0,d^{-1},t} > 0$ and $N_{0,0,d^{-1},t} > 0$. For every g such that $D_{g,t-1}^1 = 0$, $D_{g,t}^1 = 1$, and $D_{g,t}^{-1} = D_{g,t-1}^{-1} = d^{-1}$, we have

$$E(Y_{g,t} - Y_{g,t-1} | \mathbf{D}) = E(\Delta_{g,t}^1 | \mathbf{D}) + E(Y_{g,t}(0, d^{-1}) - Y_{g,t-1}(0, d^{-1}) | \mathbf{D}). \quad (28)$$

Under Assumptions 2 and 4, for all $t \geq 2$, there exists $\psi_{0,d^{-1},t} \in \mathbb{R}$ such that for all $g \in \mathcal{G}_{0,0,d^{-1},t} \cup \mathcal{G}_{1,0,d^{-1},t}$,

$$\begin{aligned} E(Y_{g,t}(0, d^{-1}) - Y_{g,t-1}(0, d^{-1}) | \mathbf{D}) &= E(Y_{g,t}(0, d^{-1}) - Y_{g,t-1}(0, d^{-1}) | \mathbf{D}_g) \\ &= E(Y_{g,t}(0, d^{-1}) - Y_{g,t-1}(0, d^{-1}) | D_{g,t-1}^1 = 0, D_{g,t-1}^{-1} = d^{-1}) \\ &= \psi_{0,d^{-1},t}. \end{aligned} \quad (29)$$

As a result,

$$\begin{aligned} & N_{1,0,d^{-1},t} E(\text{DID}_{+,d^{-1},t}^f | \mathbf{D}) \\ &= \sum_{g \in \mathcal{G}_{1,0,d^{-1},t}} N_{g,t} E(\Delta_{g,t}^1 | \mathbf{D}) + \sum_{g \in \mathcal{G}_{1,0,d^{-1},t}} N_{g,t} E(Y_{g,t}(0, d^{-1}) - Y_{g,t-1}(0, d^{-1}) | \mathbf{D}) \\ &\quad - \frac{N_{1,0,d^{-1},t}}{N_{0,0,d^{-1},t}} \sum_{g \in \mathcal{G}_{0,0,d^{-1},t}} N_{g,t} E(Y_{g,t}(0, d^{-1}) - Y_{g,t-1}(0, d^{-1}) | \mathbf{D}) \\ &= \sum_{g \in \mathcal{G}_{1,0,d^{-1},t}} N_{g,t} E(\Delta_{g,t}^1 | \mathbf{D}) + \psi_{0,d^{-1},t} \left(\sum_{g \in \mathcal{G}_{1,0,d^{-1},t}} N_{g,t} - \frac{N_{1,0,d^{-1},t}}{N_{0,0,d^{-1},t}} \sum_{g \in \mathcal{G}_{0,0,d^{-1},t}} N_{g,t} \right) \\ &= \sum_{g \in \mathcal{G}_{1,0,d^{-1},t}} N_{g,t} E(\Delta_{g,t}^1 | \mathbf{D}). \end{aligned}$$

The first equality follows by (28), the second by (29), and the third after some algebra. Given that $\text{DID}_{+,d^{-1},t}^f = 0$ if $N_{1,0,d^{-1},t} = 0$ or $N_{0,0,d^{-1},t} = 0$, we obtain, by definition of \mathcal{S}_1 and with the

convention that sums over empty sets are 0,

$$E\left(N_{1,0,d^{-1},t}\text{DID}_{+,d^{-1},t}^f \mid \mathbf{D}\right) = E\left(\sum_{\substack{g:D_{g,t}^1=1, D_{g,t}^{-1}=d^{-1} \\ (g,t) \in \mathcal{S}_1}} N_{g,t}\Delta_{g,t}^1 \mid \mathbf{D}\right). \quad (30)$$

A similar reasoning yields, for all $t \geq 2$ and $d^{-1} \in \{0, 1\}^{K-1}$,

$$E\left(N_{0,1,d^{-1},t}\text{DID}_{-,d^{-1},t}^f \mid \mathbf{D}\right) = E\left(\sum_{\substack{g:D_{g,t}^1=0, D_{g,t}^{-1}=d^{-1} \\ (g,t) \in \mathcal{S}_1}} N_{g,t}\Delta_{g,t}^1 \mid \mathbf{D}\right). \quad (31)$$

Plugging (30) and (31) into (27) yields

$$\begin{aligned} E(\text{DID}_M^f) &= E\left(E\left(\sum_{t=2}^T \sum_{d^{-1} \in \{0,1\}^{K-1}} \sum_{\substack{g:D_{g,t}^{-1}=d^{-1} \\ (g,t) \in \mathcal{S}_1}} N_{g,t}\Delta_{g,t}^1 \mid \mathbf{D}\right)\right) \\ &= E\left(E\left(\sum_{(g,t) \in \mathcal{S}_1} N_{g,t}\Delta_{g,t}^1 \mid \mathbf{D}\right)\right) \\ &= \delta_1. \end{aligned}$$

A.3 Theorem 5

As in Subsection 4.5, the proof is conditional on \mathbf{D} .

1. First, under Assumption 7, we have $\delta_1 = \sum_{s=1}^S \alpha_s \delta_{1s}$. Thus, using again Assumption 7 but also Assumption 6,

$$\text{DID}_M - \delta_1 = \sum_{s=1}^S \alpha_s (\text{DID}_s - \delta_{1s}) \sim \mathcal{N}\left(0, \sigma^2 \sum_{s=1}^S \alpha_s^2 \left(\frac{1}{n_{1s}} + \frac{1}{n_{0s}}\right)\right). \quad (32)$$

For the same reasons, we have

$$\frac{\widehat{V}}{\sigma^2} \sim \sum_{s=1}^S \alpha_s^2 \left[\frac{W_{1s}}{n_{1s}(n_{1s} - 1)} + \frac{W_{0s}}{n_{0s}(n_{0s} - 1)} \right],$$

where the $(W_{01}, W_{11}, \dots, W_{0S}, W_{1S})$ are mutually independent, independent of DID_M and $W_{ks} \sim \chi^2(n_{ks} - 1)$, by Cochran's theorem. By definition of T (see (18)), this implies that

$$\frac{\text{DID}_M - \delta_1}{\sqrt{\widehat{V}}} \sim T.$$

The result follows.

2. Without loss of generality, we assume hereafter that $G \geq G_0$, so that \mathcal{S}_G does not vary across G and has cardinal \bar{S} .

Consider the ratio $R := (\text{DID}_M - \delta_1)/\widehat{V}^{1/2}$. We first show that as $G \rightarrow \infty$, $R \xrightarrow{d} \mathcal{N}(0, 1)$. For any $(k, s) \in \{0, 1\} \times \{1, \dots, \bar{S}\}$, let $\mu_{ks} := E[\Delta Y_{g,t_s}]$ and $\sigma_{ks}^2 := V(\Delta Y_{g,t_s})$. Given that $n_{ks} \rightarrow \infty$, we have, by the central limit theorem,

$$\frac{\overline{\Delta Y}_{ks} - \mu_{ks}}{\sqrt{\sigma_{ks}^2/n_{ks}}} \xrightarrow{d} \mathcal{N}(0, 1).$$

Remark that $\text{DID}_M = \sum_{s=1}^{\bar{S}} \alpha_s (\overline{\Delta Y}_{1s} - \overline{\Delta Y}_{0s})$ and $\delta_1 = \sum_{s=1}^{\bar{S}} \alpha_s (\mu_{1s} - \mu_{0s})$. Moreover, $(\overline{\Delta Y}_{01}, \overline{\Delta Y}_{11}, \dots, \overline{\Delta Y}_{0\bar{S}}, \overline{\Delta Y}_{1\bar{S}})$ are mutually independent by Assumptions 2 and 6. Then, by, e.g., Lemma C.5 of D'Haultfœuille and Tuvaandorj (2022), we obtain

$$\frac{\text{DID}_M - \delta_1}{\sqrt{\sum_{s=1}^{\bar{S}} \alpha_s^2 (\sigma_{1s}^2/n_{1s} + \sigma_{0s}^2/n_{0s})}} \xrightarrow{d} \mathcal{N}(0, 1). \quad (33)$$

Moreover, by the law of large numbers,

$$\frac{1}{n_{1s} - 1} \sum_{g \in \mathcal{G}_{1s}} (\Delta Y_{g,t_s} - \overline{\Delta Y}_{1s})^2 \xrightarrow{P} \sigma_{ks}^2.$$

Then, $\alpha_s^2 \leq 1$ and $\min_{k,s} \liminf_G n_{ks}/G > 0$ implies that

$$G \left[\widehat{V} - \sum_{s=1}^{\bar{S}} \alpha_s^2 (\sigma_{1s}^2/n_{1s} + \sigma_{0s}^2/n_{0s}) \right] \xrightarrow{P} 0. \quad (34)$$

Next,

$$G \sum_{s=1}^{\bar{S}} \alpha_s^2 (\sigma_{1s}^2/n_{1s} + \sigma_{0s}^2/n_{0s}) \geq \left(\min_{k,s} \sigma_{ks}^2 \right) \sum_{s=1}^{\bar{S}} \alpha_s^2 \geq \frac{\min_{k,s} \sigma_{ks}^2}{\bar{S}},$$

where the latter holds by convexity of $x \mapsto x^2$ and $\sum_{s=1}^{\bar{S}} \alpha_s = 1$. Hence, by Assumption 8, we obtain

$$\liminf_G \sum_{s=1}^{\bar{S}} \alpha_s^2 (\sigma_{1s}^2/n_{1s} + \sigma_{0s}^2/n_{0s}) > 0.$$

Combined with (34), this yields

$$\frac{\widehat{V}}{\sum_{s=1}^{\bar{S}} \alpha_s^2 (\sigma_{1s}^2/n_{1s} + \sigma_{0s}^2/n_{0s})} \xrightarrow{P} 1. \quad (35)$$

Taken together, (33) and (35) imply that $R \xrightarrow{d} \mathcal{N}(0, 1)$.

Next, we prove that $T \xrightarrow{d} \mathcal{N}(0, 1)$. First, for all (k, s) , $n_{ks} \rightarrow \infty$ by Assumption 8. Thus, by the law of large numbers, $W_{ks}/(n_{ks} - 1) \xrightarrow{P} 1$. In turn, using $\liminf_G n_{ks}/G > 0$, we obtain

$$G \left[\sum_{s=1}^{\bar{S}} \alpha_s^2 \left(\frac{W_{1s}}{n_{1s}(n_{1s} - 1)} + \frac{W_{0s}}{n_{0s}(n_{0s} - 1)} \right) - \sum_{s=1}^{\bar{S}} \alpha_s^2 \left(\frac{1}{n_{1s}} + \frac{1}{n_{0s}} \right) \right] \xrightarrow{P} 0.$$

Moreover, since $G/n_{ks} \geq 1$ for all k, s , we have

$$G \sum_{s=1}^{\bar{S}} \alpha_s^2 \left(\frac{1}{n_{1s}} + \frac{1}{n_{0s}} \right) \geq 2 \sum_{s=1}^{\bar{S}} \alpha_s^2 \geq 1/\bar{S}.$$

As a result, $\liminf_G G \sum_{s=1}^{\bar{S}} \alpha_s^2 (1/n_{1s} + 1/n_{0s}) > 0$. Hence,

$$\frac{\sum_{s=1}^S \alpha_s^2 [W_{1s}/[n_{1s}(n_{1s} - 1)] + W_{0s}/[n_{0s}(n_{0s} - 1)]]}{\sum_{s=1}^S \alpha_s^2 (1/n_{1s} + 1/n_{0s})} \xrightarrow{P} 1.$$

Thus, by definition of T , $T \xrightarrow{d} \mathcal{N}(0, 1)$.

By continuity of the normal distribution, this implies that $q_{1-\alpha} \rightarrow \Phi^{-1}(1 - \alpha/2)$. Now, note that

$$P(\delta_1 \in \text{CI}_{1-\alpha}^{\text{ex}}) = F_{|R|}(q_{1-\alpha}),$$

where $F_{|R|}$ denotes the cumulative distribution function of R , which converges to $x \mapsto \max(0, 2\Phi(x) - 1)$ by what precedes. Moreover, by Pólya's theorem, the convergence is uniform. The result follows.

Web Appendix of “Two-way Fixed Effects and Differences-in-Differences Estimators with Several Treatments”

Clément de Chaisemartin Xavier D’Haultfoeille*

1 (g, t) cells belonging to \mathcal{D}_1 and not to \mathcal{S}_1

(g, t) cells belonging to \mathcal{D}_1 and not to \mathcal{S}_1 satisfy one of the five mutually exclusive conditions below.

Condition 1: (g, t) is such that $D_{g,t}^1 = 1$ for all t . As this cell’s first treatment never changes, its effect cannot be identified under a parallel trends assumption.

Condition 2: (g, t) is such that $D_{g,t}^1 = 1$, $D_{g,t-1}^1 = 1$ or $t = 1$, and $D_{g,t+1}^1 = 0$. To estimate the first-treatment’s effect in cell (g, t) , one could compare g ’s $t+1$ -to- t outcome evolution to that of another group experiencing no treatment change and with the same treatments at period $t+1$. This estimator is unbiased for the first-treatment’s effect in cell (g, t) under Assumption 5, but not under Assumption 4.

Condition 3: (g, t) is such that $D_{g,t}^1 = 1$, $D_{g,t-1}^1 = 1$ or $t = 1$, and $D_{g,t'}^1 = 0$ for a t' lower than $t-2$ or greater than $t+2$. To identify the effect of the first treatment in one such cell, one could compare g ’s t' -to- t outcome evolution to that of another group g' with the same treatments at t' and t ($D_{g',t} = D_{g',t'}$) and with the same treatments as g in t' ($D_{g,t'} = D_{g',t'}$). However, such a comparison relies on a parallel trends assumption over a longer time horizon than Assumption 4, which only imposes parallel trends over consecutive periods. One may argue that if one is ready to impose parallel trends over consecutive periods, the cost of further imposing parallel trends over longer horizons is minimal. But if one views parallel trends as a reasonable first-order approximation, rather than an assumption that exactly holds, this first-order approximation may become poorer over longer horizons, for instance if there are group specific linear trends (see Roth, Forthcoming; Rambachan and Roth, 2019). Accordingly, estimators relying on long-run parallel trends assumptions may be more biased than estimators relying on parallel trends over consecutive periods.

*de Chaisemartin: Sciences Po (email: clement.dechaisemartin@sciencespo.fr); D’Haultfoeille: CREST-ENSAE (email: xavier.dhaultfoeille@ensae.fr).

Condition 4: (g, t) is such that $D_{g,t}^1 = 1$, $D_{g,t-1}^1 = 0$, and $D_{g,t}^{-1} \neq D_{g,t-1}^{-1}$, meaning that one of g 's other treatments also changes from $t - 1$ to t . To identify the effect of the first treatment in one such cell, one could compare g 's outcome evolution to that of another group experiencing no change of its first treatment and the same changes of its other treatments as g , like in the example in Equation (8). However, such a comparison is only valid if the effect of the other treatments is constant between groups, as discussed in Section 3.2.1.

Condition 5: (g, t) is such that $D_{g,t}^1 = 1$, $D_{g,t-1}^1 = 0$, $D_{g,t}^{-1} = D_{g,t-1}^{-1}$ but there is no g' such that $D_{g',t} = D_{g',t-1} = D_{g,t-1}$, meaning that no other group experiences no treatment change and has the same baseline treatments as g . To identify the effect of the first treatment in one such cell, one could compare g 's outcome evolution to that of another group experiencing no change of its treatments and with different baseline treatments, as in the example in Equation (11). However, such a comparison is only valid if the effects of the other treatments are constant over time, as discussed in Section 3.2.1.

2 Estimators of instantaneous and dynamic effects of several treatments

In the previous subsection, we have implicitly assumed that the treatments do not have dynamic effects, since the outcome of a unit at period t only depended on her period- t treatment, not on her previous treatments.¹ When treatments can have dynamic effects, estimating the effect of a treatment controlling for other treatments is difficult. We propose an estimation strategy when there are two binary treatments, which both follow a staggered adoption design. For any $g \in \{1, \dots, G\}$, let $F_g^1 = \min\{t : D_{g,t}^1 = 1\}$ denote the first date at which group g receives the first treatment, with the convention that $F_g^1 = T + 1$ if group g never receives that treatment. Similarly, let $F_g^2 = \min\{t : D_{g,t}^2 = 1\}$ denote the first date at which group g receives the second treatment, with the convention that $F_g^2 = T + 1$ if group g never receives that treatment.

Assumption 9 (*Staggered design with two binary treatments*) For all $(g, t) \in \{1, \dots, G\} \times \{2, \dots, T\}$, $D_{g,t}^1 \in \{0, 1\}$, $D_{g,t}^2 \in \{0, 1\}$, $D_{g,t-1}^1 \leq D_{g,t}^1$, $D_{g,t-1}^2 \leq D_{g,t}^2$, and $F_g^2 \geq F_g^1$.

Assumption 9 requires that both treatments weakly increase over time, which means that once a group has switched from untreated to treated, it cannot switch back to being untreated. Assumption 9 also requires that groups start receiving the second treatment after the first. This is typically satisfied when the second treatment is a reinforcement of the first. Our running example will be that of a researcher seeking to separately estimate the effects of medical and

¹On the other hand and as discussed above, Theorems 1 and 2 do apply to dynamic effect cases, when the other treatment variables in the regression are lags of the treatment.

recreational marijuana laws in the US: so far, states have passed the former before the latter, and none of the medical and recreational laws passed since the late 1990s have been reverted. Another example where Assumption 9 holds include voter ID laws in the US, where non-strict laws are typically passed before strict ones (see Cantoni and Pons, 2021). Another example are anti-deforestation policies, where plots of lands are typically put into a concession, and then some concessions get certified (see Panlasigui et al., 2018).

To allow for dynamic effects, we need to modify our potential outcome notation. For all $(\mathbf{d}^1, \mathbf{d}^2) \in \{0, 1\}^{2T}$, let $Y_{g,t}(\mathbf{d}^1; \mathbf{d}^2)$ denote the potential outcome of group g at period t , if her two treatments from period 1 to T are equal to $\mathbf{d}^1, \mathbf{d}^2$. This dynamic potential outcome framework is similar to that in Robins (1986). It allows for the possibility that groups' outcome at time t be affected by their past and future treatments.

Our estimator relies on the following assumptions.

Assumption 10 (*No Anticipation*) For all g , for all $(\mathbf{d}^1, \mathbf{d}^2) \in \{0, 1\}^{2T}$,

$$Y_{g,t}(\mathbf{d}^1; \mathbf{d}^2) = Y_{g,t}(d_1^1, \dots, d_t^1; d_1^2, \dots, d_t^2).$$

Assumption 10 requires that a group's current outcome do not depend on her future treatments, the so-called no-anticipation hypothesis. Abbring and Van den Berg (2003) have discussed that assumption in the context of duration models, and Malani and Reif (2015), Botosaru and Gutierrez (2018), and Sun and Abraham (2021) have discussed it in the context of DID models.

For any $j \in \{1, \dots, T\}$, let $\mathbf{0}_j$ and $\mathbf{1}_j$ denote vectors of j zeros and ones, respectively. We also adopt the convention that $\mathbf{0}_0$ and $\mathbf{1}_0$ denote empty vectors. Hereafter, we refer to $Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T)$ as group g 's -treated potential outcome at period t , her outcome if she never receives either of the two treatments. Our estimators rely on the following assumption on $Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T)$.

Assumption 11 (*Independent groups, strong exogeneity, and common trends for the never-treated outcome*) For all $t \geq 2$ and $g \in \{1, \dots, G\}$, $E(Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) - Y_{g,t-1}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D})$ does not vary across g .

Assumption 11 is an adaptation of Assumptions 2-3 to the set-up we consider in this section, and where we allow for dynamic effects.

Under Assumption 9, the estimators of instantaneous and dynamic treatment effects proposed in de Chaisemartin and D'Haultfœuille (2021) can still be used with two treatments, redefining the treatment as $\tilde{D}_{g,t} = D_{g,t}^1 + D_{g,t}^2$. However, those estimators will average together effects of the first and of the second treatment. Estimating separately the effect of the first treatment is straightforward: one can just compute the estimators in Callaway and Sant'Anna (2021) or de Chaisemartin and D'Haultfœuille (2021), restricting the sample to all (g, t) s such that

$D_{g,t}^2 = 0$. In the marijuana laws example, to estimate the effect of medical marijuana laws, one can just restrict the sample to all state \times year (g, t) such that state g has not passed a recreational law yet in year t . The horizon until which dynamic effects can be estimated will just be truncated by the second treatment.

Estimating separately the effect of the second treatment is more challenging but can still be achieved, under the following, supplementary assumption.

Assumption 12 (*Restriction on the effect of the first treatment*) For all $g \in \{1, \dots, G\}$, $j \in \{1, \dots, T\}$, and $t > j$, there exists $\lambda_{j,g}(\mathbf{D})$ and $\mu_{j,t}(\mathbf{D})$ such that

$$E(Y_{g,t}((\mathbf{0}_{j-1}, \mathbf{1}_{T-(j-1)}); \mathbf{0}_T) - Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}) = \lambda_{j,g}(\mathbf{D}) + \mu_{j,t}(\mathbf{D}).$$

Assumption 12 requires that the effect of the first treatment evolves over time in the same way in every group: for any $t > j + 1$,

$$E(Y_{g,t}((\mathbf{0}_{j-1}, \mathbf{1}_{T-(j-1)}); \mathbf{0}_T) - Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}) - E(Y_{g,t-1}((\mathbf{0}_{j-1}, \mathbf{1}_{T-(j-1)}); \mathbf{0}_T) - Y_{g,t-1}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}),$$

the difference between group g 's effect of being treated for $t - (j - 1)$ and $t - 1 - (j - 1)$ periods, should be the same in every group. Still, Assumption 12 allows such treatment effects to vary in an unrestricted way with the number of time periods over which a group has been treated, and to vary with the time period at which the treatment was adopted. It also allows, to some extent, the treatment effect to vary across groups: groups' treatment effects can be arbitrarily heterogeneous at the first period where they start receiving the treatment, but the period to period evolution of that effect should then be the same in every group.

To understand why that assumption is needed, let us go back to the marijuana law example. Without Assumption 12, a state passing a recreational marijuana law may start experiencing a different outcome trend than other states that have only passed a medical law, either because of the recreational law, or because its evolution of the effect of the medical law differs from that in other states. In other words, that assumption is key to disentangle the effects of the two treatments, which is often of interest. Under the standard parallel trends assumption on the never-treated outcome (Assumption 11), one could only estimate the combined effects of the two treatments.

Though it is arguably strong, this assumption is partly testable, as we explain in more details below: it implies that groups that start receiving the treatment at the same time should then have the same outcome evolution until they adopt the second treatment. A violation of this assumption would lead our estimators to be upward (resp. downward biased) if the effect of the first treatment increases less (resp. more) in groups that adopt the second treatment than in groups that do not adopt it.

Assumption 13 (*Non-pathological design*) *There exists $(g, g') \in \{1, \dots, G\}^2$ such that $F_g^1 = F_{g'}^1$ and $1 < F_g^2 < F_{g'}^2$.*

For any $f \in \{1, \dots, T\}$, let

$$\mathcal{G}_f = \{g \in \{1, \dots, G\} : F_g^1 = f\}$$

denote the set of groups that started receiving the first treatment at date f . Let

$$\mathcal{F} = \{f \in \{1, \dots, T\} : \exists (g, g') \in \mathcal{G}_f^2 : 1 < F_g^2 < F_{g'}^2\}$$

be the set of dates such that at least two groups start receiving the first treatment at that date and start receiving the second treatment at different dates. Assumption 13 ensures that \mathcal{F} is not empty. For any $f \in \mathcal{F}$, let

$$NT_f = \max_{g \in \mathcal{G}_f} F_g^2 - 1$$

be the last period at which at least one group that started receiving the first treatment at period f has still not received the second treatment. Then, we let

$$L_{nt,f} = NT_f - \min_{g \in \mathcal{G}_f : F_g^2 \geq 2} F_g^2$$

denote the number of time periods between the first date at which a group that started receiving the first treatment at date f starts receiving the second treatment, and the last date at which a group that started receiving the first treatment at date f has not received the second treatment yet. Note that $L_{nt,f} \geq 0$ for all $f \in \mathcal{F}$. Let also

$$L_{nt} = \max_{f \in \mathcal{F}} L_{nt,f}.$$

For any $\ell \in \{0, \dots, L_{nt}\}$, $f \in \mathcal{F}$ such that $NT_f \geq \ell + f + 1$, and $t \in \{\ell + f + 1, \dots, NT_f\}$, let

$$N_{t,\ell}^f = \sum_{g \in \mathcal{G}_f : F_g^2 = t - \ell} N_{g,t}$$

denote the population at t of groups that started receiving the first treatment at date f and the second treatment ℓ periods ago at t , and such that at least one group also started receiving the first treatment at date f and has not started receiving the second treatment yet at t . Let

$$N_\ell = \sum_{f \in \mathcal{F} : NT_f \geq \ell + f + 1} \sum_{t = \ell + f + 1}^{NT_f} N_{t,\ell}^f$$

be the total population in groups reaching ℓ periods after they started receiving the second treatment at a date where there is still a group that started receiving the first treatment at the same date as their group and that has not received the second treatment yet. Across those

groups, the average cumulative effect of having received the second treatment for $\ell + 1$ periods while fixing the first treatment at its observed value is

$$\delta_\ell = E \left[\sum_{f \in \mathcal{F}: NT_f \geq \ell + f + 1} \sum_{t = \ell + f + 1}^{NT_f} \sum_{g \in \mathcal{G}_f: F_g^2 = t - \ell} \frac{N_{g,t}}{N_\ell} Y_{g,t}(\mathbf{D}_g^1; (\mathbf{0}_{t-\ell-1}, \mathbf{1}_{\ell+1})) - Y_{g,t}(\mathbf{D}_g^1; \mathbf{0}_t) \right].$$

Remark that by construction, $N_\ell > 0$ for all $\ell \in \{0, \dots, L_{nt}\}$, so δ_ℓ is well-defined for such ℓ . Note also that δ_ℓ does not include the effect of the second treatment for groups that start receiving the two treatments at the same time. For those groups, it is impossible to separately estimate the effects of the first and second treatments, using our DID estimation strategy at least.

We now define an estimator of δ_ℓ . For any $f \in \mathcal{F}$ and t such that $NT_f \geq t$, let

$$N_t^{nt,f} = \sum_{g \in \mathcal{G}_f: F_g^2 > t} N_{g,t}.$$

Then, for any $\ell \in \{0, \dots, L_{nt}\}$, $f \in \mathcal{F}$ such that $NT_f \geq \ell + f + 1$, and $t \in \{\ell + f + 1, \dots, NT_f\}$, we define

$$\text{DID}_{t,\ell}^f = \sum_{g \in \mathcal{G}_f: F_g^2 = t - \ell} \frac{N_{g,t}}{N_{t,\ell}^f} (Y_{g,t} - Y_{g,t-\ell-1}) - \sum_{g \in \mathcal{G}_f: F_g^2 > t} \frac{N_{g,t}}{N_t^{nt,f}} (Y_{g,t} - Y_{g,t-\ell-1})$$

if $N_{t,\ell}^f > 0$ and $N_t^{nt,f} > 0$, and we let $\text{DID}_{t,\ell}^f = 0$ if $N_{t,\ell}^f = 0$ or $N_t^{nt,f} = 0$. Then, for all $\ell \in \{0, \dots, L_{nt}\}$, we let

$$\text{DID}_\ell = \sum_{f \in \mathcal{F}: NT_f \geq \ell + f + 1} \sum_{t = \ell + f + 1}^{NT_f} \frac{N_{t,\ell}^f}{N_\ell} \text{DID}_{t,\ell}^f.$$

Theorem 6 *Suppose that Assumptions 1 and 9-13 hold. Then, $E[\text{DID}_\ell] = \delta_\ell$ for all $\ell \in \{0, \dots, L_{nt}\}$.*

DID_ℓ can be computed by the `did_multiplegt` Stata command, restricting the sample to the (g, t) s such that $D_{g,t}^1 = 1$, and including F_g^1 in the `trends_nonparam` option. The asymptotic normality of the DID_ℓ estimators, when the number of groups goes to infinity, could be established under similar assumptions and using similar arguments as those used to show Theorem 4 in de Chaisemartin and D’Haultfœuille (2021).

Beyond the somewhat complicated notation above, the idea underlying DID_ℓ is actually quite simple: it amounts to comparing the outcome evolution of groups that adopt/do not adopt the second treatment, and that adopted the first treatment at the same date. This ensures that the “treatment” and “control” groups involved in this comparison have been exposed to the first treatment for the same number of periods. Under Assumptions 11 and 12, this in turn ensures that their outcome evolution would have been the same if the “treatment groups” had not adopted the second treatment. The estimation procedure we propose can easily be extended

to more than two treatments. For instance, if there was a third treatment following a staggered adoption design and always adopted after the second one, one could estimate its effect using the `did_multiplt` Stata command, restricting the sample to the (g, t) s such that $D_{g,t}^2 = 1$, and including the interaction of F_g^1 and F_g^2 in the `trends_nonparam` option.

Theorem 6 complements the pioneering work of Callaway and Sant’Anna (2021) and Sun and Abraham (2021), who provide DID estimators of the effect of a single treatment following a staggered adoption design. To our knowledge, our paper is the first to consider the case with several treatments following consecutive staggered designs, which arises relatively often, as the examples given above show. Our main insight is to show that one can obtain unbiased estimators of the effect of the second treatment, under the restriction on the effect of the first treatment stated in Assumption 12, and provided one controls for the first treatment’s adoption date.

The assumptions underlying Theorem 6 are refutable. They imply that groups that start receiving the treatment at the same time should have the same outcome evolution until they adopt the second treatment, see Equation (36) in the Appendix. This can be tested, using similar placebo estimators as in de Chaisemartin and D’Haultfœuille (2021), the main difference being that one should compare groups with the same value of F_g^1 . The placebo estimators one can use to test the assumptions underlying Theorem 6 can also be computed by the `did_multiplt` command, restricting the sample to the (g, t) s such that $D_{g,t}^1 = 1$, including F_g^1 in the `trends_nonparam` option, and requesting the `placebo` option. One should still keep in mind that such pre-trends tests come with some caveats, as shown by Roth (Forthcoming): they may be underpowered and could fail to detect violations of the assumptions, and they may lead to pre-testing issues. However, our placebo estimators can be used to conduct the sensitivity analysis proposed by Manski and Pepper (2018) or Rambachan and Roth (2019).

The estimation strategy proposed in Theorem 6 requires that there is at least one pair of groups that receive the first treatment at the same date, and such that the first group receives the second treatment strictly before the second group. When the number of groups is relatively low (e.g.: the 50 US states), there may not be any pair of groups receiving the first treatment at the same time period. Then, two alternative estimation strategies can be proposed. First, instead of Assumption 12, one may assume that the effect of the first treatment evolves linearly with the number of periods of exposure, with a slope that differs across groups:

$$E(Y_{g,t}((\mathbf{0}_{j-1}, \mathbf{1}_{T-(j-1)}); \mathbf{0}_T) - Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}) = \lambda_{j,g}(\mathbf{D}) + \mu_g(\mathbf{D})(t - j).$$

Then, one can recover the counterfactual outcome that a group adopting the second treatment would have obtained without it by extrapolating a linear estimate of its outcome evolution prior to adoption. The resulting estimator can be computed by the `did_multiplt` command, restricting the sample to the (g, t) s such that $D_{g,t}^1 = 1$, and including the group indicator in the `trends_lin` option. Second, one could also strengthen Assumption 12, by assuming that

the effect of the first treatment evolves potentially non-linearly with the number of periods of exposure, but that this evolution is the same in every group and at every time period:

$$E(Y_{g,t}(\mathbf{0}_{j-1}, \mathbf{1}_{T-(j-1)}); \mathbf{0}_T) - Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}) = \lambda_{j,g}(\mathbf{D}) + \mu_{t-j}(\mathbf{D}).$$

Then, one can recover the counterfactual outcome that a group adopting the second treatment would have obtained without it, by extrapolating the outcome evolution experienced by groups that reached a similar number of periods of exposure to the first treatment without adopting the second one. The resulting estimator can be computed by the `did_multipllegt` command, restricting the sample to the (g, t) s such that $D_{g,t}^1 = 1$, and including indicators for reaching 1, 2, etc. periods of exposure to the first treatment in the `controls` option.

The approach in Theorem 6 can easily be extended to some instances where the assumptions of Theorem 6 fail. For example, there may applications with two binary treatments following a staggered adoption design, but such that some groups receive treatment 1 before treatment 2, other groups receive treatment 2 first, and other groups receive both treatments at the same time. Then, one can start by restricting attention to the subsample of groups such that $D_{g,t}^1 \geq D_{g,t}^2$ for all t and $F_g^1 < F_g^2$. This subsample includes groups that only receive treatment 1, groups that receive both treatments but receive the second one strictly after the first, and groups that do not receive any treatment. In that subsample, one can estimate the instantaneous and dynamic effects of receiving only the first treatment, using the `did_multipllegt` command and restricting the sample to (g, t) s such that $D_{g,t}^2 = 0$. One can then estimate the effect of receiving the second treatment when one has already received the first one, using the `did_multipllegt` command, restricting the sample to the (g, t) s such that $D_{g,t}^1 = 1$ and including F_g^1 in the `trends_nonparam` option. Second, one can restrict attention to the subsample of groups such that $D_{g,t}^2 \geq D_{g,t}^1$ for all t and $F_g^2 < F_g^1$. In that subsample, one can estimate the effect of receiving only the second treatment, and the effect of receiving the first treatment when one has already received the second one, using the same steps as above but reverting the roles of the first of second treatments. Finally, one can restrict attention to the subsample of groups that either receive both treatments at the same time or that do not receive any treatment, and estimate the effect of receiving both treatments at the same time using the `did_multipllegt` command. Comparing these five sets of estimates may be indicative of whether the treatments are complements or substitutes, even though differences could also be driven by heterogeneous effects across the various subsamples.

The approach outlined above can also be used when a single treatment changes several times over the duration of the panel, to isolate the effect of each change. To simplify, take the example of a binary treatment that can switch at most once from 0 to 1, and then once from 1 to 0. To estimate the effect of switching from 0 to 1, one can just use the `did_multipllegt` command in the subsample of (g, t) s such that group g has never switched from 1 to 0 at or before t . To

estimate the effect of switching from 1 to 0, one can just use the `did_multiplt` command in the subsample of (g, t) s such that group g has switched from 0 to 1 at or before t , including the date of that switch in the `trends_nonparam` option, and defining the treatment as an indicator for switching from 1 to 0. More generally, assume one is interested in the effect of a single treatment, that may not be binary and that can change multiple times, and one is interested in separately estimating the effect of each treatment change. Following similar steps as those used in the proof of Theorem 6, one can show that the estimators computed by the `did_multiplt` command, restricting the sample to the (g, t) s such that g has experienced a first treatment change at or before t and has not experienced a third treatment change at or before t , and including the date of the first treatment change interacted with groups' treatment values before and after that change in the `trends_nonparam` option, are unbiased for the instantaneous and dynamic effects of a second treatment change, under assumptions similar to Assumptions 11 and 12. One can proceed similarly to estimate effects of a third, fourth, etc. treatment change, but the corresponding estimators may soon become noisy, especially so when the treatment is not binary. In such instances, the approach proposed in de Chaisemartin and D'Haultfœuille (2021) of estimating the total cumulative effects of all treatment changes, rather than trying to separately estimate their effects, may be more feasible in practice.

3 Asymptotic inference

In this section, we show the asymptotic normality of DID_M^f and construct asymptotically valid (in fact, conservative) confidence intervals for δ_1 . The approach followed here is very similar to that considered in de Chaisemartin and D'Haultfœuille (2020) and de Chaisemartin and D'Haultfœuille (2021). Also, note that the same reasoning directly applies to δ_2 and $1/2[\delta_1 + \delta_2]$.

We use the same notation as in Section 4.5 and the proof of Theorem 5. Specifically, s refers to an observed switch, associated with a final period t_s , an initial value of $D_{g,t}^1, d_s^1$, a final value of $D_{g,t}^1, d_s^1$ and a value of $D_{g,t}^{-1}, d_s^{-1}$, constant between $t_s - 1$ and t_s . Then, define

$$U_{G,g} = \frac{G}{\sum_{s=1}^S n_{1s} |d_s^1 - d_s^{1'}|} \sum_{s=1}^S N_{g,t_s} \text{sgn}(d_s^1 - d_s^{1'}) \left[\mathbb{1}\{g \in \mathcal{G}_{1s}\} - \mathbb{1}\{g \in \mathcal{G}_{0s}\} \frac{n_{1s}}{n_{0s}} \right] \Delta Y_{g,t_s},$$

$$\hat{\sigma}^2 = \frac{1}{G} \sum_{g=1}^G \left(U_{G,g} - \text{DID}_M^f \right)^2.$$

The confidence interval of nominal level $1 - \alpha$ that we consider is

$$\text{CI}_{1-\alpha}^a = \left[\text{DID}_M^f \pm z_{1-\alpha/2} \hat{\sigma} G^{-1/2} \right],$$

where $z_{1-\alpha/2}$ is the quantile of order $1 - \alpha/2$ of the standard normal distribution.

Our results rely on Assumptions 14-15 below. Hereafter, we let $\Sigma_g := V(Y_{g,1}, \dots, Y_{g,T} | \mathbf{D})$ and let $\underline{\rho}(\Sigma)$ denote the smallest eigenvalue of any symmetric semidefinite positive matrix Σ .

Assumption 14 (*Asymptotically non-pathological design*) *The support of \mathbf{D}_g does not depend on g (and denoted by \mathcal{D}) and is finite. Let us define*

$$\mathcal{S}_g = \left\{ (t, d^1, d^{1'}, d^{-1}) : P(D_{g,t}^1 = d^1, D_{g,t-1}^1 = d^{1'}, D_{g,t}^{-1} = D_{g,t-1}^{-1} = d^{-1}) > 0 \right\}.$$

Then \mathcal{S}_g does not vary with g , $\mathcal{S}_g = \mathcal{S}$. Moreover, there exists $\underline{p} > 0$ such that for all g and $(t, d^1, d^{1'}, d^{-1}) \in \mathcal{S}$, $P(D_{g,t}^1 = d^1, D_{g,t-1}^1 = d^{1'}, D_{g,t}^{-1} = D_{g,t-1}^{-1} = d^{-1}) \geq \underline{p}$ and $P(D_{g,t}^1 = d^{1'}, D_{g,t-1}^1 = d^{1'}, D_{g,t}^{-1} = D_{g,t-1}^{-1} = d^{-1}) \geq \underline{p}$.

Assumption 15 (*Regularity conditions for asymptotic normality*) *We have, for some $c > 0$ and $\underline{\rho} > 0$,*

$$\sup_{g,t} N_{g,t} < \infty, \quad \inf_{G,g \leq G} \underline{\rho}(\Sigma_g) \geq \underline{\rho} > 0 \quad a.s. \quad \text{and} \quad \sup_{G,g \leq G, \mathbf{d} \in \mathcal{D}} E[|Y_{g,t}|^{2+c} | \mathbf{D}_g = \mathbf{d}] < \infty \quad a.s.$$

Assumption 14 basically ensures that for all switch s we observe, we have $n_{1s} \rightarrow \infty$ and $n_{0s} \rightarrow \infty$ almost surely as $G \rightarrow \infty$. Note that it always holds if groups are assumed to be identically distributed. Assumption 15 ensures that we can apply the Lyapunov central limit theorem with independent but not identically distributed variables. The second condition rules out degenerate situations where the (non-trivial) linear combinations of the $(Y_{g,1}, \dots, Y_{g,T})$ in DID_M^f would actually be constant.

Theorem 7 *Suppose that Assumptions 1-2, 4 and 14-15 hold. Then, conditional on \mathbf{D} and almost surely,*

$$\sqrt{G} \frac{\text{DID}_M^f - \delta_1}{\left(\frac{1}{G} \sum_{g: D_{g,1}=0} V(U_{G,g} | \mathbf{D}) \right)^{1/2}} \xrightarrow{d} \mathcal{N}(0, 1),$$

Moreover, we have, almost surely,

$$\liminf_{G \rightarrow \infty} \Pr \left[\delta_1 \in CI_{1-\alpha}^a | \mathbf{D} \right] \geq 1 - \alpha.$$

4 Proof of the results in the Web Appendix

Theorem 6

First, by Assumption 11, for all $t \geq 2$ there is a function of \mathbf{D} $\psi_t(\mathbf{D})$ such that

$$\psi_t(\mathbf{D}) = E(Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) - Y_{g,t-1}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}). \quad (35)$$

Then, for all $1 \leq f < t \leq T$,

$$\begin{aligned}
& E[Y_{g,t}((\mathbf{0}_{f-1}, \mathbf{1}_{T-f+1}); \mathbf{0}_T) - Y_{g,t-1}((\mathbf{0}_{f-1}, \mathbf{1}_{T-f+1}); \mathbf{0}_T) | \mathbf{D}] \\
&= E[Y_{g,t}((\mathbf{0}_{f-1}, \mathbf{1}_{T-f+1}); \mathbf{0}_T) - Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}] - E[Y_{g,t-1}((\mathbf{0}_{f-1}, \mathbf{1}_{T-f+1}); \mathbf{0}_T) - Y_{g,t-1}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}] \\
&\quad + E[Y_{g,t}(\mathbf{0}_T; \mathbf{0}_T) - Y_{g,t-1}(\mathbf{0}_T; \mathbf{0}_T) | \mathbf{D}] \\
&= \mu_{f,t}(\mathbf{D}) - \mu_{f,t-1}(\mathbf{D}) + \psi_t(\mathbf{D}); \tag{36}
\end{aligned}$$

where the second equality uses (35) and Assumption 12. Then, for any $\ell \in \{0, \dots, L_{nt}\}$, $f \in \mathcal{F}$ such that $NT_f \geq \ell + f + 1$ and $t \in \{\ell + f + 1, \dots, NT_f\}$ such that $N_{t,\ell}^f > 0$ and $N_t^{nt,f} > 0$,

$$\begin{aligned}
& E(\text{DID}_{t,\ell}^f | \mathbf{D}) \\
&= \sum_{g \in \mathcal{G}_f: F_g^2 = t - \ell} \frac{N_{g,t}}{N_{t,\ell}^f} E(Y_{g,t} - Y_{g,t-\ell-1} | \mathbf{D}) - \sum_{g \in \mathcal{G}_f: F_g^2 > t} \frac{N_{g,t}}{N_t^{nt,f}} E(Y_{g,t} - Y_{g,t-\ell-1} | \mathbf{D}) \\
&= \sum_{g \in \mathcal{G}_f: F_g^2 = t - \ell} \frac{N_{g,t}}{N_{t,\ell}^f} E(Y_{g,t}(D_g^1; (\mathbf{0}_{t-\ell-1}, \mathbf{1}_{\ell+1})) - Y_{g,t}(D_g^1; \mathbf{0}_T) | \mathbf{D}) \\
&\quad + \sum_{g \in \mathcal{G}_f: F_g^2 = t - \ell} \frac{N_{g,t}}{N_{t,\ell}^f} E(Y_{g,t}(D_g^1; \mathbf{0}_T) - Y_{g,t-\ell-1}(D_g^1; \mathbf{0}_T) | \mathbf{D}) \\
&\quad - \sum_{g \in \mathcal{G}_f: F_g^2 > t} \frac{N_{g,t}}{N_t^{nt,f}} E(Y_{g,t}(D_g^1; \mathbf{0}_T) - Y_{g,t-\ell-1}(D_g^1; \mathbf{0}_T) | \mathbf{D}) \\
&= \sum_{g \in \mathcal{G}_f: F_g^2 = t - \ell} \frac{N_{g,t}}{N_{t,\ell}^f} E(Y_{g,t}(D_g^1; (\mathbf{0}_{t-\ell-1}, \mathbf{1}_{\ell+1})) - Y_{g,t}(D_g^1; \mathbf{0}_T) | \mathbf{D}). \tag{37}
\end{aligned}$$

The first equality follows from the definition of $\text{DID}_{t,\ell}^f$, and $N_{t,\ell}^f > 0$ and $N_t^{nt,f} > 0$. The second equality follows from Assumption 10. The third equality follows from (36).

By definition of NT_f , we have $N_t^{nt,f} > 0$ for all $f \in \mathcal{F}$ and t such that $NT_f \geq t$. We adopt the convention that a sum over an empty set is equal to 0. Then, for any $\ell \in \{0, \dots, L_{nt}\}$, $f \in \mathcal{F}$ such that $NT_f \geq \ell + f + 1$ and $t \in \{\ell + f + 1, \dots, NT_f\}$, Equation (37) implies that

$$\begin{aligned}
& N_{t,\ell}^f E(\text{DID}_{t,\ell}^f | \mathbf{D}) \\
&= \sum_{g \in \mathcal{G}_f: F_g^2 = t - \ell} N_{g,t} E(Y_{g,t}(D_g^1; (\mathbf{0}_{t-\ell-1}, \mathbf{1}_{\ell+1})) - Y_{g,t}(D_g^1; \mathbf{0}_T) | \mathbf{D}).
\end{aligned}$$

We obtain the result by summing over $f \in \mathcal{F}$ and t such that $NT_f \geq \ell + f + 1$ and $t \in \{\ell + f + 1, \dots, NT_f\}$, and by the law of iterated expectations.

Theorem 7

The proof is very similar to the proof of Theorem 3 in de Chaisemartin and D'Haultfœuille (2021) so we only discuss the main steps and stress the differences here.

1. Asymptotic normality

First, remark that

$$\text{DID}_M^f = \frac{1}{G} \sum_{g=1}^G U_{G,g}.$$

Then, we establish asymptotic normality by applying Lyapunov's central limit theorem for triangular arrays. To this end, we show

$$\lim_{G \rightarrow \infty} \frac{\sum_{g=1}^G E[|U_{G,g} - E[U_{G,g} | \mathbf{D}]|^{2+c} | \mathbf{D}]}{\left(\sum_{g=1}^G V(U_{G,g} | \mathbf{D})\right)^{1+c/2}} = 0 \quad \text{a.s.} \quad (38)$$

We have $U_{G,g} = \frac{G}{\sum_{s=1}^S n_{1s} |d_s^1 - d_s^{1'}|} \sum_{t=1}^T \lambda_{g,t} Y_{g,t}$ for some $\lambda_{g,t}$. Let

$$\underline{t} = \min \left\{ t : \exists (d_1, d_1', d^{-1}) : (t, d_1, d_1', d^{-1}) \in \mathcal{S} \right\}.$$

Let s be a switch such that $t_s = \underline{t}$ and $(d_s^1, d_s^{1'}, d_s^{-1}) \in \mathcal{S}$. Then, if $g \in \mathcal{G}_{1s}$, we have $\lambda_{g, \underline{t}-1} = \text{sgn}(d_s^{1'} - d_s^1)$. Then, reasoning as for obtaining (25) in de Chaisemartin and D'Haultfœuille (2021), we get

$$\left(\frac{\sum_{s'=1}^S n_{1s'} |d_{s'}^1 - d_{s'}^{1'}|}{G} \right)^2 \sum_{g=1}^G V(U_{G,g} | \mathbf{D}) \geq \underline{\rho} \times n_{1s}. \quad (39)$$

Moreover, with the same reasoning as in de Chaisemartin and D'Haultfœuille (2021), we have

$$\left(\frac{\sum_{s=1}^S n_{1s} |d_s^1 - d_s^{1'}|}{G} \right)^2 E[|U_{G,g} - E[U_{G,g} | \mathbf{D}]|^{2+c} | \mathbf{D}] \leq C_0 \left[\max_{t=1 \dots T} |\lambda_{g,t}| \right]^{2+c}$$

for some constant $C_0 > 0$. Moreover,

$$\max_{t=1 \dots T} |\lambda_{g,t}| \leq C_1 \sup_{g,t} N_{g,t} \left[1 + \max_{s=1, \dots, S} \frac{n_{1s}}{n_{0s}} \right]. \quad (40)$$

We then obtain (38) from (39) and (40) exactly as in the proof of Theorem 3 in de Chaisemartin and D'Haultfœuille (2021).

2. Asymptotic validity of $CI_{1-\alpha}^a$

Let us define

$$\bar{\sigma}_G^2 := \frac{1}{G} \sum_{g=1}^G E[(U_{G,g} - \delta_1)^2 | \mathbf{D}].$$

The exact same reasoning as in de Chaisemartin and D'Haultfœuille (2021) leads to

$$\hat{\sigma}^2 - \bar{\sigma}_G^2 \xrightarrow{P} 0, \quad (41)$$

$$\bar{\sigma}_G^2 \geq \frac{1}{G} \sum_{g: D_{g,1}=0} V(U_{G,g} | \mathbf{D}). \quad (42)$$

Moreover, from (39), Assumptions 14-15 and the weak law of large numbers, we obtain that with probability approaching one,

$$\frac{1}{G} \sum_{g=1}^G V(U_{G,g,\ell} | \mathbf{D}) \geq C_2 \underline{\rho},$$

for some $C_2 > 0$. This ensures that $|\bar{\sigma}_G / \hat{\sigma} - 1| \xrightarrow{P} 0$. Then, write

$$G^{1/2} \frac{\text{DID}_M^f - \delta_1}{\hat{\sigma}} = \frac{\left(\frac{1}{G} \sum_{g=1}^G V(U_{G,g} | \mathbf{D}) \right)^{1/2}}{\bar{\sigma}_G} \times \left[\frac{\bar{\sigma}_G}{\hat{\sigma}} \times G^{1/2} \frac{\text{DID}_M^f - \delta_1}{\left(\frac{1}{G} \sum_{g=1}^G V(U_{G,g} | \mathbf{D}) \right)^{1/2}} \right].$$

By what precedes, the first term is smaller than one and the second converges to a standard normal variable. The result follows, with the same reasoning as in de Chaisemartin and D'Haultfœuille (2021).

References

- Abbring, Jaap H, and Gerard J Van den Berg.** 2003. “The nonparametric identification of treatment effects in duration models.” *Econometrica*, 71(5): 1491–1517.
- Botosaru, Irene, and Federico H Gutierrez.** 2018. “Difference-in-differences when the treatment status is observed in only one period.” *Journal of Applied Econometrics*, 33(1): 73–90.
- Callaway, Brantly, and Pedro H.C. Sant’Anna.** 2021. “Difference-in-Differences with Multiple Time Periods.” *Journal of Econometrics*, 225: 200–230.
- Cantoni, Enrico, and Vincent Pons.** 2021. “Strict ID laws don’t stop voters: Evidence from a US nationwide panel, 2008–2018.” *The Quarterly Journal of Economics*, 136(4): 2615–2660.
- de Chaisemartin, Clement, and Xavier D’Haultfœuille.** 2020. “Two-way fixed effects estimators with heterogeneous treatment effects.” *American Economic Review*, 110(9): 2964–96.
- de Chaisemartin, Clément, and Xavier D’Haultfœuille.** 2021. “Difference-in-Differences Estimators of Intertemporal Treatment Effects.” arXiv preprint arXiv:2007.04267.
- Malani, Anup, and Julian Reif.** 2015. “Interpreting pre-trends as anticipation: Impact on estimated treatment effects from tort reform.” *Journal of Public Economics*, 124: 1–17.
- Manski, Charles F, and John V Pepper.** 2018. “How do right-to-carry laws affect crime rates? Coping with ambiguity using bounded-variation assumptions.” *Review of Economics and Statistics*, 100(2): 232–244.
- Panlasigui, Stephanie, Jimena Rico-Straffon, Alexander Pfaff, Jennifer Swenson, and Colby Loucks.** 2018. “Impacts of certification, uncertified concessions, and protected areas on forest loss in Cameroon, 2000 to 2013.” *Biological conservation*, 227: 160–166.
- Rambachan, Ashesh, and Jonathan Roth.** 2019. “An honest approach to parallel trends.” Working paper.
- Robins, James.** 1986. “A new approach to causal inference in mortality studies with a sustained exposure period-application to control of the healthy worker survivor effect.” *Mathematical modelling*, 7(9-12): 1393–1512.
- Roth, Jonathan.** Forthcoming. “Pre-test with Caution: Event-Study Estimates after Testing for Parallel Trends.” *American Economic Review: Insights*.
- Sun, Liyang, and Sarah Abraham.** 2021. “Estimating dynamic treatment effects in event studies with heterogeneous treatment effects.” *Journal of Econometrics*, 225: 175–199.