

**TOMMASO MANFÈ**

University of Chicago

**LUCA NUNZIATA**

University of Padova and IZA

**DIFFERENCE-IN-DIFFERENCE  
DESIGN WITH REPEATED  
CROSS-SECTIONS UNDER  
COMPOSITIONAL CHANGES: A  
MONTE-CARLO EVALUATION  
OF ALTERNATIVE  
APPROACHES**

**May 2023**

**Marco Fanno Working Papers – 305**

***d*SEA**

DIPARTIMENTO DI SCIENZE  
ECONOMICHE E AZIENDALI  
'MARCO FANNO'



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

# Difference-In-Difference Design With Repeated Cross-Sections Under Compositional Changes: a Monte-Carlo Evaluation of Alternative Approaches\*

Tommaso Manfè<sup>†1</sup> and Luca Nunziata<sup>‡2,3</sup>

<sup>1</sup>University of Chicago, Booth School of Business

<sup>2</sup>University of Padua

<sup>3</sup>IZA

May 14, 2023

## Abstract

We discuss the potentially severe bias in the Difference-in-Difference (DiD) design of commonly-used methods, including the regression specification known as Two-Way-Fixed-Effects (TWFE), when researchers must invoke the conditional trend assumption but the distribution of the covariates changes over time. Building on [Abadie \(2005\)](#), we propose a Double Inverse Probability Weighting (DIPW) estimator for repeated cross-sections based on both the probability of being treated and of belonging to the post-treatment period and derive its doubly-robust version (DR-DIPW), similarly as in [Sant’Anna and Zhao \(2020\)](#). Through Monte Carlo simulations, we compare its performance with a number of methods suggested by the literature, which span from the basic TWFE (and our proposed correction) to semi-parametric estimators, including those using machine-learning first-stage estimates, following [Chernozhukov et al. \(2018\)](#). Results show that DR-DIPW outperforms the other estimators in most realistic scenarios, even if TWFE corrections provide substantial benefits. Following [Sequeira \(2016\)](#), we estimate the effect of tariff reduction on bribing behavior by analyzing trades between South Africa and Mozambique during the period 2006–2014. Contrarily to the replication by [Chang \(2020\)](#), our findings show that the effect is lower in magnitude than the one presented in the original paper.

**Keywords:** Difference-in-Difference, Monte-Carlo simulations, Semi-parametric, Machine-learning.

**JEL Codes:** C10, C13, C14, C18, C23.

---

\*This paper is a revised version of the work originally presented in Tommaso Manfè’s Master Thesis in Economics and Finance at the University of Padua entitled “Beyond regression: evaluating different semi-parametric approaches and machine learning tools in the difference-in-difference design” submitted in February 2022 and [available here](#). The coding material in its latest and intermediate versions can be downloaded from <https://github.com/tommaso-manfe>. We thank Enrico Rettore, the Master in Economics and Finance examination committee at the University of Padua, Scott Cunningham, Pedro H.C. Sant’Anna, Michael Weber, Francesco Ruggieri, Pietro Ramella, and other faculty members at the University of Chicago for their precious comments and advise. The usual disclaimer applies.

<sup>†</sup>tommaso.manfe@chicagobooth.edu

<sup>‡</sup>luca.nunziata@unipd.it

# 1 Introduction

Difference-in-Difference (DiD) is a widespread research design that estimates the causal effects of a policy treatment that affects a specific group of subjects, called the treated group, while leaving unaffected another typically comparable group, referred to as the control group. The rationale of this empirical strategy is that if treated and control groups are subject to the same time trend, the control group can be used to estimate the counterfactual potential outcome for the treatment group in the absence of treatment. The difference between this imputed potential outcome and the observed outcome for the treated in the post-treatment period is then the estimated Average Treatment Effect on the Treated (ATT).

However, the parallel trend assumption is implausible if selection into treatment depends on individual characteristics that correlate with the outcome variable. A weaker assumption consists in assuming that the parallel trend hypothesis holds after conditioning on individual observable characteristics, the so-called “conditional trend assumption”. However, even if often applied in empirical studies, the traditional TWFE is potentially biased when adding covariates in its specification.

A recent literature has focused on alternative methods for dealing with potential sources of bias when the researchers observe time-invariant covariates. In this setting, [Zeldow and Hatfield \(2019\)](#) show that, when the effect of the covariate on the outcome varies over time, the TWFE specification with covariates works only on the restrictive assumption that the mean of the covariate distribution is the same among treated and controls, which is likely to hold only in randomized experiments. In addition, it implicitly assumes homogeneous treatment effects in the covariates [Sant’Anna and Zhao \(2020\)](#), and it is subject to violations of the linear functional form. On the other hand, the other semi-parametric alternatives proposed by the literature, such as Outcome Regression (OR) [Heckman et al. \(1997\)](#), Inverse Probability Weighting (IPW) [Abadie \(2005\)](#), and Doubly-Robust DiD (DRDiD) [Sant’Anna and Zhao \(2020\)](#), despite overcoming most of the TWFE limitations, all assume that the covariates are time-invariant between the pre- and post-treatment periods. As a consequence, they may deliver biased estimates of the causal effects when these assumptions are not satisfied.

Similarly as in [Caetano et al. \(2022\)](#), we evaluate a set of novel estimators that allow the distribution of the covariates to change over time but, differently from their analysis, we derive methods suited for repeated cross-sections. The main difference from the panel data case is that the researcher does not directly observe the first difference neither in the observed outcome nor in the covariates for a given individual. As a consequence, in this setting the researcher needs not only to adjust for the covariates heterogeneity between treated and control populations but also for the heterogeneity between the pre- and post-treatment periods. To this aim, we develop a Double Inverse Probability Weighting (DIPW) scheme based on both the probability of being treated and that of belonging to the post-treatment period and derive its doubly-robust version (DR-DIPW).

In addition, in order to relax the parametric assumptions of the proposed estimators, we propose and test alternative versions of DR-DIPW that employ machine learning algorithms for the first-stage estimates, following [Chernozhukov et al. \(2018\)](#). In particular, since doubly-robust estimands satisfy the Neyman orthogonality condition needed for debiased machine learning, we construct estimators based either on lasso or random forests for the first-stage estimates. In addition, our analysis evaluates the debiased orthogonal extension of Abadie’s semiparametric DiD estimator (DMLDiD) proposed by [Chang \(2020\)](#). In general, the rationale for using machine learning is that it allows not to be restricted by a specific assumption on the parametric form for both the propensity scores and the outcome models. This is potentially beneficial since the researcher is typically unaware of the correct parametric form to assume in the analysis, resulting in model misspecification.

To corroborate our findings, we conduct a series of Monte Carlo simulations in order to investigate the finite sample properties of all semi-parametric estimators presented above. The simulations are designed to test the performance of the proposed estimators on the typical problems encountered when estimating the causal effect of interest in a DiD setting, such as trends that depend on the individual level of the covariates, heterogeneous treatment effects, and, most importantly, imbalances of the covariates between treated and controls groups and in their evolution from the pre- to the post-treatment periods, where the latter is commonly referred to as compositional changes. The different methodologies are tested across three different experimental settings in a

repeated cross-sections setup in order to provide practical guidance for practitioners.

Our analysis confirms that the commonly-used standard TWFE may be severely biased under recurrent settings. However, TWFE is not the only class of estimators that performs poorly. Both the IPW methods of [Abadie \(2005\)](#) and [Chang \(2020\)](#) are severely biased, and the same applies to the OR method developed in [Heckman et al. \(1997\)](#). Despite being approximately unbiased in most of the settings, DIPW is generally outperformed by its doubly-robust version DR-DIPW. Overall, our proposed DR-DIPW method is approximately unbiased in most of the specifications and outperforms the alternative estimators in terms of bias. In comparison to the DRDiD estimator [Sant’Anna and Zhao \(2020\)](#), its corrected weighting scheme performs better when only the two propensity scores are correctly specified or when both the propensity scores and the outcome model are both misspecified. In this last scenario, the machine learning versions of DR-DIPW have even lower bias, since they do not assume an *a priori* parametric form for the DGP. Finally, the simulations show that TWFE corrections drastically reduce the bias, but are outperformed by DR-DIPW and in general by doubly-robust estimators.

Based on our Monte Carlo simulations’ results, we propose guidelines for DiD empirical studies. In repeated cross-sections, invoking the conditional parallel trend assumption forces all the mentioned estimators to control for the level of the covariates in the post-treatment period. This exposes their estimates to the issue of bad controls, which are covariates affected by the treatment itself. For this reason, these classes of estimators retrieve the direct effect of the treatment (DAT<sub>T</sub>) on the outcome, excluding the possible effect of the covariates acting as a mediator. As a result, to retrieve the ATT, additional assumptions are required. The first option is assuming the absence of the effect of the treatment on the covariates. In this case, the DAT<sub>T</sub> and the ATT coincide. The second option is assuming that the difference in trend between treated and controls is due to a common shock, as in [Caetano et al. \(2022\)](#). In this latter case, the trend of the covariates among the controls can be used to estimate the unobserved potential level of the covariates when they are not affected by the treatment.

We finally apply this identification strategy by reproducing the analysis in [Sequeira \(2016\)](#), who investigates the effect of tariff reduction on corruption behaviors by using bribe payment data on the cargo shipments transiting from South Africa into the ports

in Mozambique. By adopting our recommended empirical strategy, we can give strong evidence against the results of the replication in [Chang \(2020\)](#), since our findings show that the effect is close and even lower in magnitude than the traditional TWFE estimation present in the original paper.

The paper is organized as follows: Section [2](#) presents the baseline features of the DiD and it illustrates TWFE and alternative semi-parametric estimators, including DIPW and DR-DIPW; Section [3](#) implements Monte Carlo simulations under different scenarios to test the performance of the various estimators; Section [4](#) provides an empirical application of the results of the simulations by analyzing the effect of tariff reduction on bribing behavior between South Africa and Mozambique during the period 2006–2014, as in [Sequeira \(2016\)](#); Section [5](#) concludes with a discussion of our most relevant findings.

## 2 Identification

### 2.1 Notation and Setup

We study the baseline case where the researcher has access to two time periods of repeated cross-sections. Define the treatment variable  $D$ , where  $d \in \{0, 1\}$ ,<sup>1</sup> as the binary indicator for whether the individual  $i$  belongs to the treated group, where the  $i$  subscript is dropped for ease of notation. Similarly, define  $T$ , where  $t \in \{0, 1\}$ , as the binary indicator that takes value zero in the pre-treatment period and one in the post-treatment period. Since the treatment is assumed to take place in between the two periods, every member of the population is untreated in the pre-treatment period. We define the potential levels of the outcome variable by using indexes that refer to the potential states of the treatment, so that  $Y_{d,t}$  denotes the outcome that would be realized for a specific value of  $d$  in period  $t$ . However, for each individual only one potential outcome is observed at each time period. At  $t = 0$ , the treatment has no effect on the potential outcomes so that  $Y_{1,0} = Y_{0,0} = Y_0$ , where we refer to the realized outcome as  $Y_t$  (i.e., not indexed by  $d$ ). At  $t = 1$ , instead,  $Y_1 = dY_{1,1} + (1 - d)Y_{0,1}$ , so we just observe the potential outcome of treatment for the treated and the potential outcome of in case of no treatment for the controls. Likewise, we denote as  $X_{d,t}$  the potential level of the covariate for the treatment group  $d$  and time  $t$ , noting again that in the pre-treatment period  $X_{1,0} = X_{0,0}$ . The object we are interested in estimating is the average effect on the treated (ATT)<sup>2</sup>, which is defined as follows:

$$ATT = E(Y_{1,1} - Y_{0,1} | D = 1)$$

which is the average difference between treated and untreated potential outcomes among the treated population. The fundamental problem of causal inference is that  $Y_{0,1}$  is not observed and thus it must be imputed.

---

<sup>1</sup>Capital letters denote random variables while small letters denote specific realizations or values of such variables.

<sup>2</sup>While usually another parameter of interest is the average treatment effect on the entire population (ATE), computing such a parameter requires additional assumptions that are unlikely to hold in this context and therefore the DiD setting usually focuses on the estimation of the ATT.

## 2.2 DiD With Time-Invariant Covariates

We initially review the existing literature on DiD with covariates, which has primarily focused on covariates that do not vary in distribution between pre- and post-treatment periods. Throughout the paper, we make the following set of assumptions.

**Assumption 1.a.** (*Sampling scheme*) *The pooled repeated cross-section data  $\{Y_i, D_i, X_i, T_i\}_{i=1}^n$  consist of iid draws from the mixture distribution*

$$P(Y \leq y, D = d, X \leq x, T = t) = t \cdot \lambda \cdot P(Y_1 \leq y, D = d, X \leq x \mid T = 1) \\ + (1 - t) \cdot (1 - \lambda) P(Y_0 \leq y, D = d, X \leq x \mid T = 0)$$

where  $(y, d, x, t) \in \mathbb{R} \times \{0, 1\} \times \mathbb{R}^k \times \{0, 1\}$ , **with the joint distribution of  $(D, X)$  being invariant to  $T$ .**

Assumption 1.a is the standard assumption among DiD that rules out compositional changes, namely a time-varying distribution of observables. Note that  $T \perp\!\!\!\perp (D, X)$  is equivalent to (i)  $X \perp\!\!\!\perp T \mid D$  and (ii)  $D \perp\!\!\!\perp T$ , i.e. (i) the observed covariates of individuals within a treatment group do not change over time, and (ii) the proportion of individuals belonging to the treatment group does not vary over time. The sampling scheme allows for each observation to be randomly chosen from either  $(Y_0, D, X)$  or  $(Y_1, D, X)$  with fixed probability  $\lambda$ . Note also that since the covariates are time-invariant, they are exogenous to the treatment. In Section 2.3, we relax Assumption 1.a to allow for compositional changes.

**Assumption 2.** (*Conditional Independence/Conditional Parallel Trend*)

$$Y_{1,1}, Y_{1,0}, Y_{0,1}, Y_{0,0} \perp\!\!\!\perp (D, T) \mid X$$

Assumption 2 requires that, conditional on the observed covariates  $X$ , the potential four potential outcomes  $Y_{d,t}$  are randomly assigned to both the treatment group and the post-treatment period. Note that this condition implies the so-called “conditional parallel trend assumption”:

$$E(Y_{0,1} - Y_{0,0} \mid X, D = 1) = E(Y_{0,1} - Y_{0,0} \mid X, D = 0)$$



which is key for the identification of causal effects in the DiD design. It implies that differences over time in the expected potential outcomes in the absence of treatment are independent of whether an individual belongs to either the treated or the control group. In other words, the conditional parallel trend assumption implies that if the treated group had not been subject to the treatment, it would have evolved, conditional on  $X$ , following the same trend observed in the control group. Therefore, the inclusion of the covariates  $X$  as controls is aimed at capturing all variables that may cause different time trends. We emphasize that this is a more robust extension of the unconditional parallel trend assumption, which claims that the parallel trend holds even when not conditioning on the covariates  $X$ . However, this latter assumption seems unlikely to hold in practice.

**Assumption 3.** (*Common Support Treatment Score*)  $P[D = 1|X] < 1 - \epsilon \quad a.s.$

for some  $\epsilon > 0$ , where we define  $p(X) \equiv P[D = 1|X]$  as the treatment score. Assumption 3 implies that it is possible to observe individuals with characteristics  $X$  among both treated and controls. In other words, the conditional probability of belonging to the treatment group given  $X$  is uniformly bounded away from one, imposing that for every value of the covariates  $X$  there is at least a small chance that the unit is not treated, and in addition, the proportion of treated units is bounded away from zero, meaning that at least a small fraction of the population is treated. The common support assumption, in contrast to the previous ones, refers to observable quantities and is therefore testable. In the case common support is not verified for all values of  $X$ , researchers usually restrict the definition of average treatment effect on the treated units where  $X$  is observable among both treated and controls.

### 2.2.1 Two-Way-Fixed Effect With Covariates

Frequently, practitioners use the following regression, usually referred to as Two-Way-Fixed Effect (TWFE) with covariates, to estimate the ATT in a DiD setting:

$$Y = \alpha + \gamma T + \beta D + \delta(T \cdot D) + X'\theta + \epsilon \quad (1)$$

where  $X = (X_1, X_2, \dots, X_p)'$  is the set of covariates with coefficients  $\theta = (\theta_1, \theta_2, \dots, \theta_p)'$ ,  $\gamma$  is the constant time effect between  $t=0$  and  $t=1$ ,  $\beta$  represents the treatment-group fixed

effect, namely the differential in the potential outcome between treated and controls in both periods  $t=0$  and  $t=1$ , and  $\delta$  represents the effect of the treatment. This specification implicitly assumes that, even when the covariates are time-invariant, three additional restrictive assumptions are verified: (i) the coefficients of the covariates do not to vary over time if the treatment is not randomized ( $X$ -specific trends), (ii) homogeneous treatment effects in  $X$ , and (iii) additive linear form of how the covariates affect the outcome.

**Case (i). ( $X$ -specific trends)** Even when the covariates  $X$  are time-invariant, their coefficients to the outcome may vary over time. We can write  $X_{1,1} = X_{1,0} \equiv X_1$  and  $X_{0,1} = X_{0,0} \equiv X_0$ . Consider for simplicity just one covariate and denote  $\theta_t$  the time-varying coefficient of  $X$ , then we have:

$$\begin{aligned} E(Y_{0,0}|X, D = 0) &= \alpha_0 + \theta_0 X_0 \\ E(Y_{0,1}|X, D = 0) &= \alpha + \gamma + \theta_1 X_0 \\ E(Y_{0,0}|X, D = 1) &= \alpha + \beta + \theta_0 X_1 \\ E(Y_{0,1}|X, D = 1) &= \alpha + \gamma + \beta + \theta_1 X_1 \end{aligned}$$

Assuming that the conditional parallel trend assumption holds, we can write:

$$\begin{aligned} E(Y_{0,1} - Y_{0,0}|X, D = 1) &= E(Y_{0,1} - Y_{0,0}|X, D = 0) \\ \alpha + \gamma + \beta + \theta_1 X_1 - (\alpha + \beta + \theta_0 X_1) &= \alpha + \gamma + \theta_1 X_0 - (\alpha + \theta_0 X_0) \\ (\theta_1 - \theta_0) \cdot (X_1 - X_0) &= 0 \end{aligned} \tag{2}$$

where the last line rearranges the terms. This implies that for covariates that do not vary over time, TWFE identifies the ATT if either: (i) the means of the covariates are the same across groups or (ii) the effects of the covariates on the outcome variable are equal in the pre and post-treatment periods ([Zeldow and Hatfield, 2019](#)).

**Case (ii). (Heterogeneous effects)** In most realistic settings, the effect of the treatment is likely to vary for different values of the covariates  $X$ . However, TWFE and its correction implicitly assume homogeneous treatment effects in  $X$  and therefore, when this additional restriction is not satisfied, the estimated causal parameter may differ

from the true ATT (Meyer, 1995; Abadie, 2005; Sant’Anna and Zhao, 2020; Roth et al., 2022). For instance, let the treatment effect be heterogeneous in  $X$ , as in Cunningham (2021), namely redefining the potential outcomes for the treated in the post period as  $E(Y_{1,1}|X, D = 1) = \alpha + \gamma + \beta + (\delta + \rho X_1) + \theta X_1$ . Then, even assuming time-invariant coefficients of the covariates  $\theta_1 = \theta_0$  we have:

$$\begin{aligned} E(Y_{1,1} - Y_{1,0}|X, D = 1) &= \delta + [(\theta + \rho)X_1 - \theta X_1] \\ &= \delta + \rho X_1 \end{aligned}$$

while the TWFE estimate of the ATT is just  $\delta$ . As a consequence, whenever  $\rho \neq 0$  and thus the treatment is heterogeneous in  $X$ , the regression estimate does not identify the true ATT, even when the covariates are restricted to be time-invariant.

**Case (iii). (Non-additive linear form of the CEF for the covariates)** Since in most settings it is not possible to use a fully saturated model in  $X$ , TWFE assumes a CEF that is a linear function of  $X$ , so that the regression equation might differ from the true CEF. Indeed, the linear specification for the control variables implies that the assumption of common trends is conditional on the linear index  $X'\theta$  which is more restrictive than assuming common trends conditional on  $X$ . For example, if the vector  $X$  does not affect the potential outcome linearly, then the potential outcome is:

$$\begin{aligned} E(Y_{d,t}|X) &= f(\alpha + \gamma T + \beta D + \delta TD + \theta X) \\ &\neq \alpha + \gamma T + \beta D + \delta TD + \theta X \end{aligned}$$

and the TWFE estimate is biased since the model does not capture non-linearities.

The standard TWFE specification can be improved by allowing some corrections. Zeldow and Hatfield (2019) argue that by adding an interaction between the time dummy  $T$  and the time-invariant covariates  $X$ , the confounder effect of the covariates in presence of homogeneous treatment effects in  $X$  can be eliminated. A natural extension that allows for both  $X$ -specific trends and heterogenous effects can be written as:

$$Y = \alpha + \gamma T + \beta D + \delta(TD) + X'_i\theta + (TX')\omega + (DX')\nu + (TDX')\rho + \epsilon \quad (3)$$

where the interaction term  $T \cdot D \cdot X'$  explicitly allows for the effect of the treatment to change depending on the level of  $X$ .

In addition to the suggested TWFE correction, the literature has focused on time-invariant controls methods which rely less on parametric assumptions and may have other desirable properties, like double robustness. In the next sections, we review the properties of these estimators.

### 2.2.2 Outcome Regression

The outcome regression (OR henceforth) approach relies on the researchers' ability to correctly specify a model for the evolution of the outcome of interest. Intuitively, since the ATT under conditional parallel trends requires the computation of four conditional expectations of the outcome variable, OR computes the quantities that are not directly observable by specifying a functional model based on the covariates  $X$ . More precisely, the two conditional expectations of the observed outcome in the pre- and post- treatments periods for the treated are directly computed by means of sample averages. The remaining two conditional expectations are predicted for the treated sub-population on the basis of estimations obtained on the controls. In other words, an outcome model is estimated on untreated units given their covariate values, and then fitted values are predicted using the empirical distribution of  $X$  among treated units. Usually, the outcome model is estimated through regression, hence the alternative name of regression adjustment (RA), but other more flexible non-parametric methods can be employed as well, such as nearest neighbor matching, which associates treated with untreated units with close covariate values.

More formally, following Heckman et al. (1997), starting from the definition of the ATT under conditional parallel trends and using the law of iterated expectations, we obtain:

$$\begin{aligned} ATT &= E[E(Y_1 - Y_0|X, D = 1) - E(Y_1 - Y_0|X, D = 0)|D = 1] \\ &= E(Y_1 - Y_0|D = 1) - E[E(Y_1 - Y_0|X, D = 0)|D = 1] \end{aligned} \quad (4)$$

where the first term in Eq. (4) can be computed by taking sample averages, while the second expected value must be estimated. One way of estimating it is by fitting a regression

on the controls group data and taking predictions based on the empirical distribution of  $X$  among treated units. More formally:

$$\delta^{OR} = \bar{Y}_{1,1} - \bar{Y}_{1,0} - \left[ \frac{1}{n_{treat}} \sum_{i|D_i=1} (\hat{\mu}_{0,1}(X_i) - \hat{\mu}_{0,0}(X_i)) \right] \quad (5)$$

where  $\bar{Y}_{d,t} = \sum_{i|D_i=1} Y_{it}/n_{d,t}$  is the sample average outcome among treated units in treatment group  $d$  at time  $t$ , and  $\hat{\mu}_{d,t}(X)$  is an estimator of the true, unknown  $m_{d,t}(x) \equiv E[Y_t|D = d, X = x]$ , which is usually estimated by running a regression on data from the observed control sub-population defined by  $d$  and  $t$  and deriving fitted values from the empirical distribution of  $X$  among the treated individuals. Intuitively, when using a linear specification for  $\hat{\mu}_{d,t}(X)$ , the model would be close to the version of TWFE with covariates that includes also all the interactions between  $X_i$  and both treatment group and time dummies, as in Eq. (3). The two models differ because the outcome regression approach adopts a re-weighting scheme based on the distribution of  $X$  among units with  $D = 1$  (Roth et al., 2022). The condition for the consistency of the ATT of the outcome regression is the correct specification of  $\hat{\mu}_{d,t}(X)$ .

### 2.2.3 Inverse Probability Weighting

The Inverse Probability Weighting (IPW) approach proposed by Abadie (2005) avoids the direct modelling the outcome evolution. Its focus is on the treatment model, namely the conditional probability of being in the treatment given the set of covariates,  $p(X) \equiv P(D = 1|X)$ . The idea of the IPW estimator is to adjust for confounding factors using the propensity score to balance baseline individual characteristics in the treated and untreated groups. Since we are dealing with repeated cross-sections, it is not directly possible to apply the usual IPW weighting scheme to the first difference in outcomes within individuals, as it happens in the case of panel data. As a result, the following

weighting scheme is suggested:

$$\begin{aligned}
E[Y_{11}] &= E \left[ \frac{YDT}{\lambda p(X)} \right] \\
E[Y_{10}] &= E \left[ \frac{YD(1-T)}{(1-\lambda)p(X)} \right] \\
E[Y_{01}] &= E \left[ \frac{Y(1-D)T}{\lambda(1-p(X))} \right] \\
E[Y_{00}] &= E \left[ \frac{Y(1-D)(1-T)}{(1-\lambda)(1-p(X))} \right]
\end{aligned}$$

where  $\lambda = E[T]$ , i.e. the proportion of individuals observed in  $t = 1$ . To see why, consider for example:

$$\begin{aligned}
E \left[ \frac{YDT}{\lambda p(X)} \right] &= E \left[ E \left[ \frac{YDT}{\lambda p(X)} \middle| X \right] \right] \\
&= E \left[ E \left[ \frac{YT}{\lambda p(X)} \middle| X \right] p(D = 1|X) \right] \\
&= E \left[ E \left[ \frac{Y}{\lambda p(X)} \middle| D = 1, X \right] p(T = 1|D, X) p(D = 1|X) \right] \\
&= E \left[ E \left[ \frac{Y\lambda}{\lambda p(X)} \middle| T = 1, D = 1, X \right] p(X) \right] \\
&= E [E [Y_{11}|X]] \\
&= E [Y_{11}]
\end{aligned}$$

where we exploited the law of iterated expectations, the notion that  $p(X) = P(D = 1|X)$  is known given  $X$ , and that  $E(YDT) = 1 \cdot E(YT|D)P(D = 1) + 0 \cdot E(YT|D)P(D = 0)$ . A similar reasoning applies to the other three remaining potential outcomes. By taking the difference in differences of the potential outcomes, namely  $E[(Y_{11} - Y_{10}) - (Y_{01} - Y_{00})] \equiv E[\Delta Y_1 - \Delta Y_0]$ , we obtain:

$$\begin{aligned}
E[\Delta Y_1 - \Delta Y_0] &= E \left[ \frac{YDT - YD\lambda}{\lambda(1-\lambda)p(X)} - \frac{YT - YDT - Y\lambda + YD\lambda}{\lambda(1-\lambda)(1-p(X))} \right] \\
&= E \left[ Y \cdot \frac{(T - \lambda)}{\lambda(1-\lambda)} \frac{(D - p(X))}{(1-p(X))p(X)} \right]
\end{aligned}$$

which coincides with the ATE. The ATT is instead retrieved with the following estimand:

$$\delta^{IPW} = \frac{1}{E(D) \cdot \lambda} \cdot E \left[ \frac{D - p(X)}{1 - p(X)} \cdot \frac{T - \lambda}{1 - \lambda} \cdot Y \right] \quad (6)$$

which can be estimated using the following sample analog:

$$\delta^{IPW} = \frac{1}{\lambda \cdot \frac{1}{n} \sum_{j=1}^n (D_j)} \cdot \sum_{i=1}^n \left[ \frac{D_i - \hat{\pi}(X_i)}{1 - \hat{\pi}(X_i)} \cdot \frac{T_i - \lambda}{1 - \lambda} \cdot Y_i \right] \quad (7)$$

Intuitively, IPW produces a weighting scheme that weights-down the observed outcome  $Y_{d,t}$  for the individuals with covariate values over-represented among their time and treatment group category (namely, with low  $\frac{p(X)}{1-p(X)}$ ), and weights-up the observed outcome for the individuals with covariate values under-represented among their group. Consequently, the adjustment balances the distribution of covariates between treated and untreated groups. The unknown propensity score  $p(X) = P(D = 1|X)$  is usually estimated by means of logistic regression or a linear probability model, even if non-parametric models can be employed as well. The IPW approach will generally be consistent when the propensity score model is correctly specified.

#### 2.2.4 Doubly Robust Difference-in-Difference

[Sant'Anna and Zhao \(2020\)](#) combine the OR and the IPW approaches into a doubly robust estimand for the ATT. The double robustness property means that if either the propensity score model or the outcome regression models are misspecified (but not both), the resulting estimand still identifies the ATT. Intuitively, the Doubly Robust Difference-in-Difference (DRDiD) estimator they propose has the advantages of each of the two individual DiD methods and, at the same time, circumvents some of their weaknesses. Denote  $\mu_{d,t}(X)$  as the arbitrary model for the true, unknown conditional expectation function  $m_{d,t}(x) \equiv E[Y|D = d, T = t, X = x]$ , and for ease of notation define  $\mu_{d,Y}(T, X) \equiv T \cdot \mu_{d,1}(X) + (1 - T) \cdot \mu_{d,0}(X)$ , where recall  $d, t \in \{0, 1\}$ . Intuitively,  $\mu_{d,Y}(T, X)$  represents the outcome model for a given treatment group  $D = d$ . Then the

estimand is defined as:

$$\delta_1^{dr} = E \left[ \left( \omega_1(D, T) - \omega_0(D, T, X; p) \right) \left( Y - \mu_{0,Y}(T, X) \right) \right] \quad (8)$$

where:

$$\begin{aligned} \omega_1(D, T) &= \omega_{1,1}(D, T) - \omega_{1,0}(D, T) \\ \omega_0(D, T, X; p) &= \omega_{0,1}(D, T, X; p) - \omega_{0,0}(D, T, X; p) \end{aligned}$$

and for  $t \in \{0, 1\}$ :

$$\begin{aligned} \omega_{1,t}(D, T) &= \frac{D \cdot 1\{T = t\}}{E[D \cdot 1\{T = t\}]} \\ \omega_{0,t}(D, T, X; p) &= \frac{(1 - D)p(X) \cdot 1\{T = t\}}{1 - p(X)} \bigg/ E \left[ \frac{(1 - D)p(X) \cdot 1\{T = t\}}{1 - p(X)} \right] \end{aligned}$$

The relative sample analog is obtained by replacing  $p(x)$  with  $\hat{\pi}$  and the expectation with sample means. The first term of  $\delta_1^{dr}$  represents the IPW weighting scheme based on the propensity score, while the second term represents the outcome regression part of the estimand. [Sant'Anna and Zhao \(2020\)](#) present also a locally semi-parametrically efficient version of the above estimator, which is characterized by an asymptotic variance that achieves the semi-parametric efficiency bound when the propensity score and outcome regression are correctly specified:

$$\begin{aligned} \delta_2^{dr} &= \delta_1^{dr} + (E[\mu_{1,1}(X) - \mu_{0,1}(X)|D = 1] - E[\mu_{1,1}(X) - \mu_{0,1}(X)|D = 1, T = 1]) \\ &\quad - (E[\mu_{1,0}(X) - \mu_{0,0}(X)|D = 1] - E[\mu_{1,0}(X) - \mu_{0,0}(X)|D = 1, T = 0]) \end{aligned} \quad (9)$$

In the Monte Carlo simulations in [Section 3](#), we consider only the estimand  $\delta_2^{dr}$ . The outcome equation and the propensity score can be modeled either parametrically, for instance with a linear and logistic regression respectively, or non-parametrically. In the first case, the authors name the estimator DRDiD. The authors also use the inverse probability tilting estimator ([Graham et al., 2012](#)) for the treatment model and weighted least-squares for the outcome model. In this case, they name the estimator Improved



DRDiD (IMP DRDiD). We will stick to this definition in the remaining sections. The two estimators will generally be consistent if either the propensity score or the outcome model is correctly specified.

In Section 2.3, we propose a set of novel semi-parametric extensions of the models discussed above that allow the distribution of the covariates to change over time.

## 2.3 DiD Under Compositional Changes

Despite researchers more often observe time-varying covariates in empirical settings, just a few papers consider this setup. This is problematic since time-varying covariates can act as a confounder in DiD settings. For example, [Zeldow and Hatfield \(2019\)](#) shows that by allowing time-varying covariates in the commonly used TWFE regression with covariates, Eq. (2) can be rewritten:

$$E(Y_{0,1} - Y_{0,0}|X, D = 1) = E(Y_{0,1} - Y_{0,0}|X, D = 0)$$

$$\theta_1(X_{1,1} - X_{0,1}) - \theta_0(X_{1,0} - X_{0,0}) = 0$$

where we adopt the potential covariate notation  $X_{d,t}$ . Therefore, TWFE retrieves the ATT only if (i) the relationship between the covariates and the outcome is constant in time, and (ii) the difference in the mean of the covariates between treated and controls is the same in pre and post-treatment periods (i.e.  $X_{1,1} - X_{0,1} = X_{1,0} - X_{0,0}$ ) ([Zeldow and Hatfield, 2019](#)). Such a condition is particularly restrictive compared to Eq. (2) and implies that a time-varying covariate is a confounder if its relationship with the outcome is time-varying or the covariate evolves differently between the treated and control groups. Moreover, all methods considered above, i.e. OR, IPW, and DRDiD, assume time-invariant covariates and therefore might be biased from confounding effects when instead the distribution of the covariates varies over time.

Among the few papers focussing on compositional changes, [Caetano et al. \(2022\)](#) propose a doubly robust estimand for panel data which controls both for the pre- and post-treatment levels of the covariate and outline specific assumption under which the ATT is identified. Instead, [Hong \(2013\)](#) proposes a two-variate matching estimator for repeated cross-sections based on the probability of being treated in the pre- and post-treatment

periods, which are defined separately. We adopt a similar approach by specifying two propensity scores: one for the probability of being treated (treatment score), and the other for the probability of belonging to the post-treatment period (time score). We then construct an inverse probability weighting scheme based on both scores and derive its doubly-robust estimand. From now on, we rely on the following additional set of assumptions.

**Assumption 1.b.** (*Sampling scheme*) *The pooled repeated cross-section data  $\{Y_i, D_i, X_i, T_i\}_{i=1}^n$  consist of iid draws from the mixture distribution*

$$P(Y \leq y, D = d, X \leq x, T = t) = t \cdot \lambda \cdot P(Y_1 \leq y, D = d, X \leq x \mid T = 1) \\ + (1 - t) \cdot (1 - \lambda) P(Y_0 \leq y, D = d, X \leq x \mid T = 0)$$

where  $(y, d, x, t) \in \mathbb{R} \times \{0, 1\} \times \mathbb{R}^k \times \{0, 1\}$ , **with the joint distribution of  $(D, X)$  being time-varying with respect to  $T$ .**

Assumption 1.b replaces Assumption 1.a allowing covariates to be time-varying. Since now the covariates evolve over time, we cannot rule out the possibility that they are affected by the treatment.

**Assumption 4.** (*Common Support Time Score*)  $P[T = 1 \mid D, X] < 1 - \epsilon \quad a.s.$

for some  $\epsilon > 0$ . The time score  $t(D, X) \equiv P[T = 1 \mid D, X]$  is the probability of belonging to the post-treatment period conditional on the covariate  $X$  and the treatment group  $D$ . Intuitively, it allows for heterogeneous time-trends in the covariates among treated and untreated individuals. As for the treatment score, Assumption 4 ensures that, for any values of  $X$ , there will be some units in the post-treatment period for both treated and untreated units. Under time-varying covariates, it is necessary for identification purposes to have assumptions on how the covariates are allowed to vary over time. Similarly as in [Caetano et al. \(2022\)](#), we provide conditions on  $X$  that allows to identify the ATT:

**Assumption 5.a.** (*Covariates exogeneity*)  $X_{1,1} = X_{0,1} \quad \forall i : D_i = 1$

Assumption 5.a states that in the post-treatment period the potential covariate level in case of treatment is equal to the potential covariate level in case of no treatment. This

condition rules out bad controls, namely covariates affected by the treatment. Despite this may be a too restrictive assumption, in repeated cross-sections it is not possible to control for the pre-treatment levels of the covariates since each individual is observed in only one time period. Therefore, there is an inherent trade-off between correcting for the heterogeneity in the evolution of the covariates and allowing bad controls. However, an alternative assumption to retrieve the ATT in case of time-varying covariates is:

**Assumption 5.b.** (*Covariates Parallel Trend*)

$$X_{0,1} = X_{0,0} + E(X_{0,1} - X_{0,0} | D = 0) \quad \forall i : D_i = 1$$

where we follow the potential covariate level notation. Assumption 5.b imposes that the individual trend for the treated observation  $i$  would have been the same as the mean trend observed in the untreated population. Intuitively, it assumes that the differential trend in the covariates among treated and untreated units is fully caused by the treatment. The imputed  $X_{0,1}$  among treated is then replaced to the observed level of  $X_{1,1}$  and the ATT can be estimated through the methods presented in this section. Despite being a restrictive assumption, such a condition is very different from 5.a and may be used as an additional robustness check to test how the estimates are sensitive to the potential role of covariates acting as mediators.

### 2.3.1 Double inverse-probability weighting (DIPW)

Building on Abadie (2005), we propose a weighting scheme that corrects for the heterogeneous trends in  $X$  between treated and controls. Intuitively, the idea is that of replacing  $\lambda$ , which is the proportion of people observed at  $T = 1$ , by the time score  $t(D, X)$ , which is the probability of being observed at  $T = 1$  conditional on covariates  $X$  and treatment status  $D$ . The weights to retrieve the four average potential outcomes can be specified

as follows:

$$\begin{aligned}
E[Y_{1,1}] &= E \left[ \frac{YDT}{t(D, X)p(X)} \right] \\
E[Y_{1,0}] &= E \left[ \frac{YD(1-T)}{(1-t(D, X))p(X)} \right] \\
E[Y_{0,1}] &= E \left[ \frac{Y(1-D)T}{t(D, X)(1-p(X))} \right] \\
E[Y_{0,0}] &= E \left[ \frac{Y(1-D)(1-T)}{(1-t(D, X))(1-p(X))} \right]
\end{aligned}$$

More formally, by exploiting the law of iterated expectations, the conditional independence assumption 2, and basic conditional expectations rules, it is possible to show that the weighting scheme retrieves the average potential outcomes  $E[Y_{d,t}]$ . For example, in the case of  $E[Y_{1,1}]$ :

$$\begin{aligned}
E \left[ \frac{YDT}{t(D, X)p(X)} \right] &= E \left[ E \left[ \frac{YDT}{t(D, X)p(X)} \middle| X \right] \right] \\
&= E \left[ E \left[ \frac{YT}{t(D, X)p(X)} \middle| X \right] p(D = 1|X) \right] \\
&= E \left[ E \left[ \frac{Y}{t(D, X)p(X)} \middle| D = 1, X \right] p(T = 1|D = 1, X)p(D = 1|X) \right] \\
&= E \left[ E \left[ \frac{Y}{t(D, X)p(X)} \middle| T = 1, D = 1, X \right] t(D, X)p(X) \right] \\
&= E [E [Y_{1,1}|X]] \\
&= E [Y_{1,1}]
\end{aligned}$$

and similar passages applies to  $E[Y_{1,0}]$ ,  $E[Y_{0,1}]$ , and  $E[Y_{0,0}]$ . The estimator for ATE is therefore:

$$ATE = E \left[ Y \cdot \frac{(T - t(D, X))}{t(D, X)(1 - t(D, X))} \frac{(D - p(X))}{(1 - p(X))p(X)} \right] = E [Y\omega]$$

where we defined the general weighting scheme as  $\omega = \frac{(T-t(D,X))}{t(D,X)(1-t(D,X))} \frac{(D-p(X))}{(1-p(X))p(X)}$ . Since we are mainly focused on the ATT, this can be retrieved by:

$$\begin{aligned}
ATT &= E[Y\omega|D = 1] \\
&= \int E[Y\omega|D = 1, X] p(X|D = 1) \\
&= \int E[Y\omega|D = 1, X] \frac{p(X)P(D = 1|X)}{P(D = 1)} \\
&= \int E[Y\omega|D = 1, X] \frac{p(X)P(D = 1|X)}{P(D = 1)} \\
&= E[Y\omega|D = 1, X] \frac{P(D = 1|X)}{P(D = 1)} \\
&= E\left[Y\omega \frac{p(X)}{E(D)}\right]
\end{aligned}$$

The terms  $t(D, X)$  and  $p(X)$  must be estimated in the sample, usually by logit or probit regression. If the models are correctly specified, the estimator retrieves the ATT. In practice, these Horvitz-Thompson weights can take values at the extremes of  $[0, 1]$ , therefore standardizing the weights is beneficial for reducing the estimator variance. Therefore, we propose the following estimand for the ATT that uses standardized Hayek weights as follows:

$$\delta^{dipw} = E\left[\left(\omega_1^{mod}(D, T) - \omega_0^{mod}(D, T, X; p)\right)Y\right] \quad (10)$$

where:

$$\begin{aligned}
\omega_1^{mod}(D, T) &= \omega_{1,1}^{mod}(D, T) - \omega_{1,0}^{mod}(D, T) \\
\omega_0^{mod}(D, T, X; p, t) &= \omega_{0,1}^{mod}(D, T, X; p) - \omega_{0,0}^{mod}(D, T, X; p)
\end{aligned}$$

and for  $t \in \{0, 1\}$ :

$$\omega_{1,1}^{mod}(D, T; t) = \frac{D \cdot T}{t(D, X)} \Big/ E \left[ \frac{D \cdot T}{t(D, X)} \right]$$

$$\omega_{1,0}^{mod}(D, T; t) = \frac{D \cdot (1 - T)}{(1 - t(D, X))} \Big/ E \left[ \frac{D \cdot (1 - T)}{(1 - t(D, X))} \right]$$

$$\omega_{0,1}^{mod}(D, T, X; p, t) = \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot t(D, X)} \Big/ E \left[ \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot t(D, X)} \right]$$

$$\omega_{0,0}^{mod}(D, T, X; p, t) = \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot (1 - t(D, X))} \Big/ E \left[ \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot (1 - t(D, X))} \right]$$

where  $\omega_{d,t}^{mod}(D, T, X; p, t)$  are the standardized Hayek weights after some rearrangements. Note that we used the  $\omega_{d,t}^{mod}$  notation to allow for easy comparison to the weights in [Sant'Anna and Zhao \(2020\)](#).

### 2.3.2 Doubly Robust Inverse Probability Weighting (DR-DIPW)

Similarly to [Sant'Anna and Zhao \(2020\)](#), we combine the OR and the DIPW approaches into a doubly robust estimand for the ATT. As in Section 2.2.4, denote  $\mu_{d,t}(X)$  as the arbitrary model for the true, unknown conditional expectation function  $m_{d,t}(x) \equiv E[Y|D = d, T = t, X = x]$ , and for ease of notation define  $\mu_{d,Y}(T, X) \equiv T \cdot \mu_{d,1}(X) + (1 - T) \cdot \mu_{d,0}(X)$ , where  $d, t \in \{0, 1\}$ . Intuitively,  $\mu_{d,Y}(T, X)$  represents the outcome model for a given treatment group  $D = d$ . Then the new estimand is defined as:

$$\delta_1^{drdipw} = E \left[ \left( \omega_1^{mod}(D, T) - \omega_0^{mod}(D, T, X; p) \right) \left( Y - \mu_{0,Y}(T, X) \right) \right] \quad (11)$$

where:

$$\begin{aligned} \omega_1^{mod}(D, T) &= \omega_{1,1}^{mod}(D, T) - \omega_{1,0}^{mod}(D, T) \\ \omega_0^{mod}(D, T, X; p, t) &= \omega_{0,1}^{mod}(D, T, X; p) - \omega_{0,0}^{mod}(D, T, X; p) \end{aligned}$$

and for  $t \in \{0, 1\}$ :

$$\omega_{1,1}^{mod}(D, T) = \frac{D \cdot T}{t(D, X)} \Big/ E \left[ \frac{D \cdot T}{t(D, X)} \right]$$

$$\omega_{1,0}^{mod}(D, T) = \frac{D \cdot (1 - T)}{(1 - t(D, X))} \Big/ E \left[ \frac{D \cdot (1 - T)}{(1 - t(D, X))} \right]$$

$$\omega_{0,1}^{mod}(D, T, X; p, t) = \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot t(D, X)} \Big/ E \left[ \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot t(D, X)} \right]$$

$$\omega_{0,0}^{mod}(D, T, X; p, t) = \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot (1 - t(D, X))} \Big/ E \left[ \frac{(1 - D) \cdot T \cdot p(X)}{(1 - p(X)) \cdot (1 - t(D, X))} \right]$$

Again, the first term of  $\delta_1^{drdipw}$  represents the IPW weighting scheme, while the second term represents the outcome regression part of the estimand. We also present the locally semi-parametrically efficient version of the above estimator closely following [Sant'Anna and Zhao \(2020\)](#). Recall that this version is characterized by an asymptotic variance that achieves the semi-parametric efficiency bound when the working models for the nuisance functions are correctly specified:

$$\begin{aligned} \delta_2^{drdipw} = & \delta_1^{drdipw} + (E[\mu_{1,1}(X) - \mu_{0,1}(X)|D = 1] - E[\mu_{1,1}(X) - \mu_{0,1}(X)|D = 1, T = 1]) \\ & - (E[\mu_{1,0}(X) - \mu_{0,0}(X)|D = 1] - E[\mu_{1,0}(X) - \mu_{0,0}(X)|D = 1, T = 0]) \end{aligned} \quad (12)$$

The relative sample analog is obtained by replacing  $p(X)$ ,  $t(D, X)$ , and  $\mu_{d,t}(X)$  with their in-sample estimations. In the Monte Carlo simulations in [Section 3](#), we consider only the estimand  $\delta_2^{drdipw}$ . To maintain consistency and comparability with the DRDiD and IMP DRDiD estimator, we name DR-DIPW the version that employs linear and logistic regression for the outcome equation and the propensity scores respectively, while name Improved DR-DIPW (IMP DR-DIPW) the version using the inverse probability tilting estimator ([Graham et al., 2012](#)) for the treatment and time scores models and weighted least-squares for the outcome model. The two estimators will generally be consistent if either the propensity score or the outcome model is correctly specified.

### 2.3.3 Machine Learning DR-DIPW

We propose an alternative version of the DR-DIPW estimator that employs machine learning algorithms for first-stage estimates. The advantage of machine learning over traditional estimation methods is that they do not need a parametric assumption of the functional form of the model under study, avoiding the risk of misspecification. [Chernozhukov et al. \(2018\)](#) and the related literature that followed identified three main conditions that enable the use of first-stage machine learning without creating bias in the estimates of the causal parameter. The first is the so-called Neyman orthogonality condition which guarantees that the estimand must be insensitive to small perturbations of the nuisance functions. This property is typically satisfied for estimands that are doubly robust ([Chernozhukov et al., 2018](#); [Sant’Anna and Zhao, 2020](#); [Farrell et al., 2021](#)), like the DR-DIPW estimator in Eq. (11) and (12). Under this condition, machine learning estimates of the nuisance functions (in our case the propensity scores and outcome model functions) are allowed even if they are generally biased due to regularization. Indeed, using a Neyman-orthogonal score eliminates the biases arising from the first-stage estimates. The second condition refers to the rate of convergence of the machine learning estimators used for the nuisance parameters. In particular, they have to converge to the true parameter using the  $L^2(P)$  norm at a rate faster than  $o(N^{-1/4})$ . [Chernozhukov et al. \(2018\)](#) shows that such a condition is generally met by most machine learning estimators such as lasso, ridge, random forests, neural nets, and various hybrids and ensembles of these methods.

Finally, the authors suggest using a form of sample splitting: the nuisance parameters are estimated on a random partition, while the remaining sample is used for the estimation of the orthogonal score. However, we do not employ any form of sample splitting since we do not find any benefit in terms of reduction of bias when applying this technique to the doubly robust estimators in our simulation setting. Indeed, as noted in [Farrell et al. \(2021\)](#), sample splitting can also not be applied when the parameter of interest is not itself learned from the data. In the case of the DR-DIPW estimator, only the regression functions and propensity scores must be estimated and the ATT is not itself learned from the data. In what follows, we discuss and evaluate two DRDiD machine learning



extensions of the DRDiD model. The first employs lasso in both the estimation of the propensity score and the outcome model, while the second utilizes random forest for both nuisance parameters. The performance of these two estimators is compared to the traditional parametric first-stage estimates in the Monte Carlo Simulations in Section 3.

### 3 Monte Carlo Simulations

In this section, we conduct a series of Monte Carlo simulations in order to investigate the finite sample properties of the proposed estimators in a repeated cross-sections setup. The different methodologies are tested across two different experimental settings. Each design is characterized by two repeated cross-sections, one observed at  $t = 0$  and another at  $t = 1$ , with a total sample size of  $n = 1000$  observations. The Monte Carlo simulation consists of 10000 randomly generated datasets and estimation results are stored at each repetition.

Our two simulations focus on cases in which the treatment is not randomized, since in that instance all methods presented in the paper yield unbiased estimates of the ATT. Conversely, Experiment 1 allows instead for non-randomized selection into treatment and  $X$ -specific trends, while still assuming time-invariant covariates and homogeneous treatment effects. In Experiment 2, though, the distribution of covariates is allowed to vary between the pre- and post-treatment periods and the treatment effects are heterogeneous in  $X$ . For this reason, Experiment 2 reproduces the most realistic and indicative version of the data generating process (DGP, henceforth). The choice of the functional forms of our DGPs, presented below, is aimed at preserving the comparability with the work of [Sant'Anna and Zhao \(2020\)](#) and [Kang and Schafer \(2007\)](#), which employed the same functional specifications. Indeed, Experiment 1 closely reproduces the Monte Carlo simulations in [Sant'Anna and Zhao \(2020\)](#), while Experiment 2 extends the study into the aforementioned more realistic conditions.

First of all, for a generic variable  $W = (W_1, W_2, W_3, W_4)'$ , we define the underlying

true outcome and propensity score model as:

$$f_{reg}(W) = 210 + 25.4 \cdot W_1 + 13.7 \cdot (W_2 + W_3 + W_4) \quad (13)$$

$$f_{ps}(W) = 0.75 \cdot (-W_1 + 0.5 \cdot W_2 - 0.25 \cdot -0.1 \cdot W_4) \quad (14)$$

The function  $f_{ps}(W)$ , which determines selection into treatment, is modelled through the inverse of the logit function, i.e.  $expit(f_{ps}(W)) = \frac{\exp(f_{ps}(W))}{1 + \exp(f_{ps}(W))}$ , which has the desirable property of producing an average propensity score of 0.5. In other words, assuming parametrically a logit model for the propensity score (with all the relevant covariates) will lead to a correct estimation of the probability of being treated, by construction.

In the context of each of our experiments, the baseline function for the outcome  $f_{reg}(W)$  produces a mean of  $E(Y) = E[f_{reg}(W)] = 210.0$  and, when combined with  $f_{ps}(W)$ , leads to  $E(Y|D = 0) = 200.0$  and  $E(Y|D = 1) = 220.0$ . As outlined in [Kang and Schafer \(2007\)](#), the selection bias in this DGP is not severe because the difference between the average outcome among the treated units and the average outcome among the full population is only a one-quarter of a population standard deviation. Nevertheless, this difference is large enough to invalidate the performance of naive estimators.

For each of our two experiments, we replicate both scenarios when the researcher correctly specifies or misspecifies the parametric form of the model. The aim is to assess the estimators' performance in terms of bias when the researcher cannot correctly specify the functional form of the model under study. Overall, each of the two experiments considers four different DGPs.

We define the generic vector  $W$ , which can either represent vector  $Z$ , which is the set of variables observed by the researcher, or vector  $X$ , which is not observable. The idea is that the unique generic DGP (expressed in terms of  $W$ ) leads, for each experiment, to four cases depending on whether  $W$  is replaced by the observed vector  $Z$  or by the unobservable vector  $X$ . The misspecification of the models derive from the fact that  $Z$  is a highly non-linear transformation of  $X$  and its interactions. When the modelling functions are defined as  $f_{ps}(Z)$  and  $f_{reg}(Z)$ , i.e. they are functions of the observed  $Z$ , then both the propensity score and outcome regression models will be correctly specified since the variables we observe coincide with those affecting the outcome. We call this scenario DGP

A. When the data are generated by  $f_{ps}(X)$  and  $f_{reg}(X)$ , then the researcher, who has only access to  $Z$ , will misspecify both models. We call this scenario DGP D. Typically, such a scenario is the most realistic since researchers do not have an *a priori* knowledge of the phenomenon under analysis. We also consider the two cases in which just one of the two models is correctly specified. We call this scenarios DGP B (when the outcome model is correctly specified) and DGP C (when only the propensity score model is correctly specified).

More formally, assume  $X = (X_1, X_2, X_3, X_4)'$  is distributed as  $N(0, I_4)$  with  $I_4$  representing the  $4 \times 4$  identity matrix. For  $j = 1, 2, 3, 4$  define the standardized variable  $Z_j = (\tilde{Z}_j - E[\tilde{Z}_j]) / \sqrt{Var(\tilde{Z}_j)}$  where  $\tilde{Z}_1 = \exp(0.5X_1)$ ,  $\tilde{Z}_2 = 10 + X_2 / (1 + \exp(X_1))$ ,  $\tilde{Z}_3 = (0.6 + X_1X_2/25)^3$ , and  $\tilde{Z}_4 = (20 + X_2 + X_4)^2$ . Note that the non-linear transformations that generates the relationship between the individual variables of  $Z$  and  $X$  include a wide range of functional forms, such as quadratic, cubic and exponential. In addition, such transformations include interactions of the  $X$ s in order to achieve additional complexity in the relationship that links  $Z$  and  $X$ . As a result, when the propensity score and outcome regression models are misspecified, i.e. when the DGPs are built from  $X$  whereas we observe  $Z$  only, this is likely to cause a bias in the estimation. In this case, allowing for non-parametric first-stage estimates, which may better capture the non-linearities between  $Z$  and  $X$ , can minimize the bias. For example, random forest will build a non-parametric relation between  $Z$  and the outcome, approximating more closely the relation between the true  $X$  (which is a non-linear transformation of  $Z$ ) and the outcome rather than simply assuming a linear relationship that may not hold in the data. Similarly, lasso allows for more flexibility with respect to traditional methods.

Table 1 summarizes the different estimation methods that are tested in each experiment. They are evaluated in terms of average bias, root mean square error (RMSE), variance, and computational time required for the estimation. When not otherwise specified, all estimators employ a logit model for the propensity score and a linear regression model for the outcome. Therefore, the first is estimated using maximum likelihood and the second by ordinary least squares. Note that the choice of using a logit model for the propensity score is required to perfectly match, by construction, the functional form of the probability of being treated in our DGPs. When the DGPs are built from  $f_{ps}(Z)$  and

$f_{reg}(Z)$ , the models are therefore correctly specified.<sup>3</sup>

For the DR-DIPW we also allow for the possibility of first-stage estimates using machine learning methods. Such non-parametric methods should better capture the non-linearities under investigation when the working models are misspecified. When lasso is used, the outcome and the treatment model are designed as a penalized linear and a penalized logistic regression, respectively. Lasso is performed in R using the 'glmnet' package (Friedman et al., 2010) and the shrinkage parameter  $\lambda$  is selected through 10-fold cross-validation and represents the largest value of  $\lambda$  whose cross-validation error is within 1 standard deviation from the minimum. This allows to select the sparsest model with performances approximately equal to the optimum.

Since lasso implicitly performs variable selection and can therefore handle a large set of covariates, in the simulations we allow lasso to employ an expanded set of covariates that include all third order terms and interactions of the original variables. As a result, lasso performs a selection from a wider set of variables and may more precisely capture the non-linearities in the functional forms related to the phenomenon under study. On the contrary, in real data analysis, when employing traditional estimators, the researcher may be constrained by the risk of including a number of predictors  $p$  that is too large. In some extreme cases this number may be close or even higher than the number of observations  $n$ , invalidating the estimation. Despite being technically possible in our synthetic dataset to include all interactions terms for the traditional estimators as well (since we have just four regressors, the expanded set of covariates  $p = 34$  would be still lower than  $n$ ), we limit the inclusion of these higher order terms and interactions only to the lasso model to emulate the more recurrent real-world scenario where it is not possible for traditional estimators to do so.

When random forest is used, the estimation is implemented using the 'randomForest' R package (Liaw and Wiener, 2002). The number of trees is set to 500 in order to obtain a good balance between accuracy and computational effort. At each node, as common practice, the number of randomly sampled input variables is restricted to  $\sqrt{p}$ , where  $p$  is again the number of predictors (James et al., 2013).

---

<sup>3</sup>This would not be true if, for example, a linear probability model would be used for the propensity score.

When using machine learning tools, we do not perform sample splitting, even if generally suggested in [Chernozhukov et al. \(2018\)](#) and [Bach et al. \(2021\)](#). In fact, as noted in [Farrell et al. \(2021\)](#), sample splitting should be used in cases where the parameter of interest itself is learned from the data. For the causal estimands in the DR-DIPW setting, the regression functions and propensity score must be estimated, but these are first-stage estimates. The second stage of the DR-DIPW is just the expected value of the estimand of the ATT. Therefore, in this case sample splitting does not provide any advantage. Finally, it is also important noting that, when using sample splitting, the required computational time is increased by a factor of  $k$  in case of  $k$ -fold sample splitting. If the standard errors are calculated by bootstrap, this would lead to an increase of computational time in the order of  $k \times n$  where  $n$  is the number of the repetitions in the bootstrap.

In what follows we discuss each experiment under different assumptions for the DPG.

### 3.1 Experiment 1: $X$ -specific Trends and Non-Randomized Selection

Experiment 1 closely replicates the simulation presented by [Sant’Anna and Zhao \(2020\)](#). In addition to  $X$ -specific trends, here the selection into treatment is not randomized, causing additional obstacles to the identification of the causal parameter. In a non-randomized experiment setting, selection into treatment may be associated with some individual characteristics  $X$  and is therefore likely to cause heterogeneity in the distribution of the covariates between treated and controls. In addition, in Experiment 1 covariates are assumed to be time-invariant, so that compositional changes in the independent variables are ruled out, and treatment effects are homogeneous in  $X$ .

The four different DGPs are specified as in [Table 2](#), where  $\epsilon_0(d)$ ,  $\epsilon_1(d)$ ,  $d = 0, 1$ , are independent standard normal random variables representing the stochastic error term of the potential outcomes, the propensity score  $p(W)$  is a logistic transformation of the generic function  $f_{ps}(W)$ ,  $\lambda$  is the proportion of observations in  $t = 1$  and  $U_d$  and  $U_t$  are independent standard uniform stochastic variables used to randomly select individuals into treatment and post-treatment period, respectively. For a generic variable  $W$ ,  $v(W, D)$

is an independent normal random variable with mean  $D \cdot f_{reg}(W)$  and unit variance which represents the time-invariant unobserved group heterogeneity between treated and untreated populations. The trend is specified as  $\tau(W) = f_{reg}(W)$ , and therefore in the post-treatment period  $t = 1$  it sums to the standard function of the outcome model  $f_{reg}$ . This explains the presence of the factor 2 that multiplies the term  $f_{reg}(W)$  in the formula of the potential outcome  $Y_{1,1}$ . The observed outcome is  $Y = DTY_{1,1} + D(1 - T)Y_{1,0} + (1 - D)TY_{0,1} + (1 - D)(1 - T)Y_{0,0}$ . In the aforementioned DGPs, the true ATT is zero.

Figure 1 shows the distribution of the observed covariate  $X_4$  between the two groups in the pre- and post-treatment periods (plots with similar characteristics can be shown for the remaining three regressors  $X_1$ ,  $X_2$ , and  $X_3$ ). Within each treatment group category there are no significant changes between  $t = 0$  and  $t = 1$  because the covariates are time-invariant. However, the plot displays heterogeneity in the distribution of the covariates among treated and controls. Note that this heterogeneity is not severe, since the difference between the mean of the treated and the mean of the full population is constructed to be only one-quarter of a population standard deviation. If an estimator performs poorly in such a scenario, it may even be more biased when treated and control groups are more heterogeneous in terms of covariate distribution. Therefore, if traditional estimators fail under a moderate selection bias like this one, such a result strengthens the need of adopting alternative more flexible estimation techniques. The results of the simulation can be found in Tables 3 to 6.

In Experiment 1, the bias of the TWFE estimator with covariates is evident. Independently of whether the model is correctly specified, TWFE is typically severely biased, with a bias of 20.762 even in the most favorable scenario embodied by DGP A. The TWFE correction is characterized by a significantly lower bias: in Experiments 1A and 1B, where the outcome model is correctly specified, it is approximately unbiased, while in Experiments 1C and 1D, it is outperformed by most of the other estimators. The most efficient class of estimators is represented by the doubly robust methods, which are approximately unbiased when either the propensity score or the outcome model are correctly specified. They also have better in-sample properties when both models are misspecified. In the latter case indeed, IMP DRDiD and IMP DR-DIPW have approximately half of the bias (2.550 and 2.563 respectively) compared to the TWFE correction (5.108). In

Experiment 1D, the lasso version of the DR-DIPW has the lowest bias (1.938), thanks to its more flexible parametric assumptions. Conversely, the other machine learning method DMLDiD is severely biased in all four scenarios and is characterized by a very high variance in its estimates, even when IPW and DIPW have a low bias, not far from the best performing-estimators.

### 3.2 Experiment 2: $X$ -specific Trends and Non-Randomized Selection under Compositional Changes

Experiment 2 tests the proposed estimators in presence of  $X$ -specific trends, non-randomized selection into treatment, heterogeneous treatment effects in  $X$ , and compositional changes in the distribution of  $X$  between the pre and post-treatment periods. As discussed in Section 2.3, the inclusion of time-varying covariates in the TWFE is likely to yield biased estimates, and the other aforementioned alternative semi-parametric estimators may perform poorly as well since they assume time-invariant  $X$ .

Table 7 describes the four DGPs in Experiment 2. In this design, the DGPs are subject to two main changes. The first is that the covariates between treated and controls are subject to different time trends. We model this scenario by specifying the conditional probability of belonging to the post-treatment period  $t(D, X)$ . In particular, the probability is calculated as the logistic transformation of  $f_{ps}(-W)$  for treated and  $f_{ps}(W)$  for controls. This way, the observations are assigned to a specific time period so that different evolutions of the distribution of the covariates among treated and controls are produced. Figure 2 clearly elucidates the implications of the new DGPs in terms of the covariate distribution among the different groups. We can see that the distribution of  $X_4$ , in addition to being heterogeneous among treated and controls, it evolves differently for the two groups between  $t = 0$  and  $t = 1$ . The second change consists in allowing the treatment effect to vary with  $X$ . This is achieved by denoting the treatment effect as  $\tilde{\delta}(W) = -10W_1 + 10W_2 - 10W_3 - 10W_4$ . In addition, to guarantee that the ATT is zero as in the two previous experiments, we use the demeaned transformation of  $\tilde{\delta}(W)$ , e.g.  $\delta(W) = \tilde{\delta}(W) - E_{i|D=1}[\tilde{\delta}(W)]$ , where  $E_{i|D=1}[\tilde{\delta}(W)]$  denotes the ATT before demeaning. The results of the simulations are displayed in Tables 8 to 11.

Experiment 2 represents the most realistic setting for a typical researcher, and therefore its implications are particularly relevant. In most DGPs, the traditional TWFE specification is severely biased. An exception is Experiment 2D, where TWFE has the lowest bias (3.437), but this is likely caused by different sources of bias offsetting each other since in all other scenarios standard regression works poorly. However, when adding the relevant interaction terms with the covariates, the TWFE correction substantially improves the estimates. Indeed, in Experiments 2A and 2B, TWFE CORR is approximately unbiased as the doubly-robust estimators. Indeed, in the first experiments also DRDiD and IMP DR-DiD are approximately unbiased, even if they are originally built for time-invariant covariates. This indicates that their flexible specification of the outcome model can be naturally extended to time-varying covariates under our assumptions. This is not always the case, since the OR approach and the DRDiD versions that are not locally-efficient (not reported, available upon request) are substantially biased also in the case of correctly specified models. The intuition is that, under time-varying covariates, it is not sufficient to train an outcome model solely on the untreated population, and a model trained on the treated units must be specified as well. However, in all four simulations, the traditional weighting scheme of the IPW, which is also embedded in [Sant’Anna and Zhao \(2020\)](#), is markedly biased. On the contrary, the DIPW correctly accounts for compositional changes and shows very limited bias when the treatment and time scores models are correctly specified. Indeed, it is characterized by one of the best in-sample performances in terms of bias also in Experiment 2D.

Overall, the doubly robust versions of the DIPW are the models with better properties in all four simulations. They are approximately unbiased in all settings where at least either the propensity scores or the outcome models are correctly specified. In particular, they strictly outperform their estimators in [Sant’Anna and Zhao \(2020\)](#) that do not account for compositional effects in Experiment 2C and 2D, where the weighting scheme plays a more important role. In particular, DR-DIPW and DRDiD have a bias of 0.277 and 4.379, respectively, in Experiment 2C, while IMP DR-DIPW and IMP DRDiD have 0.519 and 1.092. Similarly, in Experiment 2D the bias for DR-DIPW and DRDiD is 12.793 and 17.736, respectively, and for IMP DR-DIPW and IMP DRDiD is 10.429 and 12.793. Noticeably, the LASSO DR-DIPW has the lowest bias (7.586) in Experiment



2D, which replicates the most realistic scenario where the researcher cannot know the functional form of the phenomenon under study. For this reason, we suggest using this version in empirical studies.

## 4 Empirical illustration: the effect of tariff reduction on corruption behaviors

We illustrate the implications of using alternative estimation methods by reproducing the analysis in [Sequeira \(2016\)](#) who investigate the effect of tariff reduction on corruption behaviors by using bribe payment data on the cargo shipments transiting from South Africa into the ports in Mozambique. This contribution adds to a rich debate on whether a decrease in tariff rates disincentives corruption. On the one side, tariff rates decreases are expected to lower the incidence of bribing behavior since they reduce the marginal advantage to evade taxes ([Allingham and Sandmo, 1972](#); [Poterba, 1987](#); [Fisman and Wei, 2001](#)). On the other side, lower tariff levels have also an income effect, increasing private agents' resources to pay higher bribes ([Slemrod and Yitzhaki, 2002](#); [Feinstein, 1991](#)).

In 1996, a trade agreement between South Africa and Mozambique paced a series of tariff reductions that took place between 2001 and 2015, with the largest of them occurring in 2008 and entailing an average nominal tariff rate of about 5 percentage points. In this context, [Sequeira \(2016\)](#) collected primary data on the bribe payments of shipments imported from South Africa to Mozambique from 2007 to 2013 through an audit study. As previously documented in [Sequeira and Djankov \(2014\)](#), it was common for cargo owners, in exchange for tariff evasion, or simply to avoid the threat of being cited for real or fictitious irregularities, to bribe border officials in charge of collecting all tariff payment and of providing clearance documentation. For example, prior to 2008, approximately 80 percent of the random sample of tracked shipments were linked to sizeable bribe payments during the clearing process (mean bribes reached USD 128 per tonnage). As a consequence, [Sequeira \(2016\)](#) exploits the exogenous change in tariffs induced by the trade agreement to examine the effect of changes in tariffs on corruption levels. Since not all products experienced a variation in tariff rates during this period, the

author adopts a Difference-in-Difference design to isolate the causal relationship between tariffs and corruption, on pooled cross sectional data collected between 2007 and 2013, for a total of 1084 observations. More specifically, the design is based on the canonical TWFE estimator in the following specification:

$$\begin{aligned}
y_{it} = & \gamma_1(TariffChangeCategory_i \times POST) + \mu POST \\
& + \beta_1 TariffChangeCategory_i + \beta_2 BaselineTariff_i \\
& + \Gamma_i + p_i + \omega_t + \delta_i + \epsilon_{it}
\end{aligned} \tag{15}$$

where  $y_{it}$  represents the natural log of the amount of bribe paid for shipment  $i$  in period  $t$ , conditional on paying a bribe,  $TariffChangeCategory_i \in \{0, 1\}$  takes value one if the commodity was subject to tariff reduction,  $POST \in \{0, 1\}$  denotes the years following 2008, and  $BaselineTariff_i$  is a control for the pre-treatment tariff for product  $i$ . The specification also accounts for a vector of product, shipment, clearing agent, and firm-level characteristics  $\Gamma_i$  which includes the elements summarised in Table 12. Industry, year, and clearing agent fixed effects are included, denoted by  $p_i$ ,  $\omega_t$ , and  $\delta_i$  respectively. The parameter of interest is the coefficient of the interaction between the time and treatment dummies, namely  $\gamma_1$ .

The main finding of Sequeira (2016) is that the tariff reduction led to a significant drop in the amount of bribe paid. Chang (2020) replicates the estimation by using the DMLDiD estimator. In Table 13, we report the results obtained by the two authors, where TWFE refers to the standard specification in Sequeira (2016) (Equation 1 of Table 9 in their paper), TWFE ( $\Gamma_i \times POST$ ) is the specification that also includes the interactions between the covariates  $\Gamma_i$  and  $POST$  (which differs from Eq. (3) where all the interactions between the covariates and the time and treatment group dummies are added), while DMLDiD is estimated by either using kernel or lasso in the first-stage estimates. Overall, the DMLDiD estimates claim that the effect of the reduction was larger than originally thought.

However, both classes of estimators used in this analysis are likely to be characterized by a substantial degree of bias. Figure 3 shows the standardized mean difference in trend for each of the 33 covariates between treated and controls. The latter is defined

as  $\frac{(\bar{X}_{11}-\bar{X}_{10})-(\bar{X}_{01}-\bar{X}_{00})}{std(X)}$ , where the overbar indicates the mean. The graph suggests that the distribution of the covariates is time-varying because, in the case of time-invariant  $X$ , the metrics should be 0, and the presence of heterogeneous covariate trends between treated and controls. This is the condition tested in Experiment 2 in Section 3, where we assume a non-randomized treatment scenario with  $X$ -specific trends, compositional changes, heterogeneous effects, and potential non-linearities in the DGP. Because in our simulation the two estimators were strongly biased, computing the effect of tariff reduction on bribing patterns with TWFE and DMLDiD may be misleading. In addition, DMLDiD estimates in Table 13 suffer from very high standard errors which blur the interpretation of the empirical findings. For example, the 95 percent confidence interval lies approximately between 0.318 and  $-14.306$  for the kernel DMLDiD, and the same applies to its lasso version, even if to a smaller degree.

Motivated by these considerations, we employ the lasso DR-DIPW estimator, which proved to be the least biased estimator in the most realistic setting of the Monte Carlo simulations (Experiment 2D), to the current study of the effect of tariff reduction on bribing behaviour. In such a setting, the low number of observations also does not allow for traditional first-stage estimation methods to produce accurate fitted values, favouring the use of sparse machine learning methods like lasso. Indeed, the sample is characterized by a large number of observations for the control group in the post-treatment period, but a limited number of observations for the other three groups, namely the treated in the pre and post-treatment period (120 and 56 observations respectively), and the controls in the pre-treatment period (84).

The lasso specification captures non-linearities by allowing for a richer set of covariates:  $\Gamma_i$  is expanded to include all second order terms and interactions, leading to a set of 112 controls. Contrarily to Chang (2020), we stick to the specification in Sequeira (2016) by including all industry, time and clearing agent fixed-effects, even if their interactions are not generated due to computational tractability. The DR-DIPW standard errors are computed through weighted bootstrap, similarly to Sant’Anna and Zhao (2020).

The ATT is estimated first invoking Assumption 5.a, where we rule out the presence of bad controls. We then adopt Assumption 5.b for a subset of covariates, assuming that the evolution of these covariates among treated and controls is due to a common shock

and the variation of the covariates among the treated in the post-treatment period is due to the effect of the treatment. More precisely, we calculate the imputed covariate level for a subset of covariates that we deem potentially affected by the treatment, like the log value of the shipment per tonnage and the industry classification category of the products shipped. Indeed, the reduction in tariff may alter the typology of products traded and therefore these covariates may act as bad controls in our estimates. Despite both assumptions being stringent, they define very different scenarios and thus they can be used as a valuable robustness check.

Our final estimates are displayed in Table 14. Our results, across different methods and specifications, corroborate the hypothesis that the tariff reduction led to a drop in the amount of bribe paid, but give compelling evidence against the assumption that the effect was higher in magnitude. In fact, the standard TWFE seems to overestimate the ATT ( $-3.748$ ), while the lasso DR-DIPW estimates are  $-2.478$  and  $-1.594$  under Assumptions 5.a and 5.b, respectively. The standard errors are typically lower than those in Chang (2020), producing more precise confidence intervals estimations.

In summary, our findings reveal that the tariff reduction had a significant effect on bribing behavior, but the impact is smaller than originally estimated by the standard TWFE specification and by DMLDiD as in Chang (2020). The DATT is  $-2.478$ , while when we assume that some covariates are affected by the treatment, the estimate is even lower in magnitude, namely  $-1.594$ .

## 5 Conclusion

Our analysis shows that the commonly-used DiD estimators, including TWFE, may be severely biased when invoking the conditional parallel trend assumption and the covariates vary over time. We specify a set of inverse probability weights (DIPW) that use both a treatment and a time score to retrieve the DATT under compositional changes and, building on [Sant'Anna and Zhao \(2020\)](#), we propose its doubly-robust version (DR-DIPW). When comparing the performance of the various estimators proposed by the literature in Monte Carlo simulations, we show that DR-DIPW tends to outperform all other available methods. In particular, when we assume that the researcher is not able to correctly specify the functional form of the propensity score and outcome models, the machine learning version of DR-DIPW with lasso first-stage estimates has even a relatively more limited bias with respect to traditional first-stage estimates.

We then propose two different assumptions to investigate the possible role of the covariates as mediators of the treatment effect. The first and more common one is to assume the absence of bad controls, namely that the treatment does not affect the covariates. In this case, the DATT and the ATT coincide. Alternatively, it is possible to assume that the covariates are subject to a common shock between the pre- and post-treatment period, among treated and controls. In this case, it is possible to impute the level of the potential covariate level of the treated by assuming a common trend with the observed variation among the control group.

We apply both identification strategies by reproducing the analysis in [Sequeira \(2016\)](#) which investigates the effect of tariff reduction on corruption behaviors by using bribe payment data on the cargo shipments transiting from South Africa into the ports in Mozambique. Our estimates show that tariff reduction led to a decrease in bribes paid but the effect is significantly lower in magnitude than the one estimated using the TWFE specification with covariates in the original paper and a DMLDiD specification as in the replication by [Chang \(2020\)](#).

## References

- Alberto Abadie. Semiparametric difference-in-differences estimators. *The Review of Economic Studies*, 72(1):1–19, 2005.
- Michael G. Allingham and Agnar Sandmo. Income tax evasion: a theoretical analysis. *Journal of Public Economics*, 1(3-4):323–338, 1972. URL <https://EconPapers.repec.org/RePEc:eee:pubeco:v:1:y:1972:i:3-4:p:323-338>.
- Philipp Bach, Victor Chernozhukov, Malte S Kurz, and Martin Spindler. Doubleml— an object-oriented implementation of double machine learning in r. *arXiv preprint arXiv:2103.09603*, 2021.
- Carolina Caetano, Brantly Callaway, Stroud Payne, and Hugo Sant’Anna Rodrigues. Difference in differences with time-varying covariates, 2022.
- Neng-Chieh Chang. Double/debiased machine learning for difference-in-differences models. *The Econometrics Journal*, 23(2):177–191, 2020.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 01 2018. ISSN 1368-4221. doi: 10.1111/ectj.12097. URL <https://doi.org/10.1111/ectj.12097>.
- Scott Cunningham. A tale of time varying covariates. <https://causalinf.substack.com/p/a-tale-of-time-varying-covariates>, 2021. Accessed: 07/02/2022.
- Max H Farrell, Tengyuan Liang, and Sanjog Misra. Deep neural networks for estimation and inference. *Econometrica*, 89(1):181–213, 2021.
- Jonathan S. Feinstein. An econometric analysis of income tax evasion and its detection. *RAND Journal of Economics*, 22(1):14–35, 1991. URL <https://EconPapers.repec.org/RePEc:rje:randje:v:22:y:1991:i:spring:p:14-35>.
- Raymond Fisman and Shang-Jin Wei. Tax Rates and Tax Evasion: Evidence from “Missing Imports” in China. NBER Working Papers 8551, National Bu-

- reau of Economic Research, Inc, October 2001. URL <https://ideas.repec.org/p/nbr/nberwo/8551.html>.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1): 1–22, 2010. URL <https://www.jstatsoft.org/v33/i01/>.
- Bryan S. Graham, Cristine Campos De Xavier Pinto, and Daniel Egel. Inverse Probability Tilting for Moment Condition Models with Missing Data. *The Review of Economic Studies*, 79(3):1053–1079, 04 2012. ISSN 0034-6527. doi: 10.1093/restud/rdr047. URL <https://doi.org/10.1093/restud/rdr047>.
- James J Heckman, Hidehiko Ichimura, and Petra E Todd. Matching as an econometric evaluation estimator: Evidence from evaluating a job training programme. *The review of economic studies*, 64(4):605–654, 1997.
- Seung-Hyun Hong. Measuring the effect of napster on recorded music sales: difference-in-differences estimates under compositional changes. *Journal of Applied Econometrics*, 28(2):297–324, 2013.
- Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning*, volume 112. Springer, 2013.
- Joseph DY Kang and Joseph L Schafer. Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data. *Statistical science*, 22(4):523–539, 2007.
- Andy Liaw and Matthew Wiener. Classification and regression by randomforest. *R News*, 2(3):18–22, 2002. URL <https://CRAN.R-project.org/doc/Rnews/>.
- Breed D Meyer. Natural and quasi-experiments in economics. *Journal of business & economic statistics*, 13(2):151–161, 1995.
- James M Poterba. Tax Evasion and Capital Gains Taxation. *American Economic Review*, 77(2):234–239, May 1987. URL <https://ideas.repec.org/a/aea/aecrev/v77y1987i2p234-39.html>.

Jonathan Roth, Pedro HC Sant'Anna, Alyssa Bilinski, and John Poe. What's trending in difference-in-differences? a synthesis of the recent econometrics literature. *arXiv preprint arXiv:2201.01194*, 2022.

Pedro HC Sant'Anna and Jun Zhao. Doubly robust difference-in-differences estimators. *Journal of Econometrics*, 219(1):101–122, 2020.

Sandra Sequeira. Corruption, trade costs, and gains from tariff liberalization: Evidence from southern africa. *American Economic Review*, 106(10):3029–63, 2016.

Sandra Sequeira and Simeon Djankov. Corruption and firm behavior: Evidence from african ports. *Journal of International Economics*, 94(2):277–294, 2014.

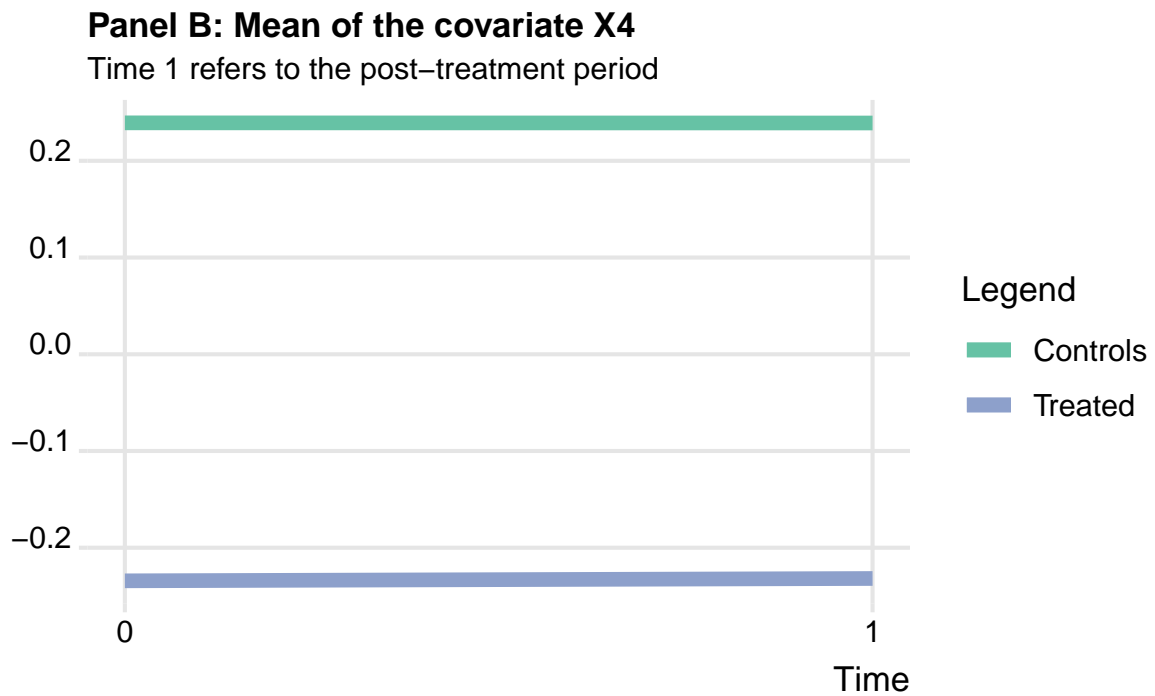
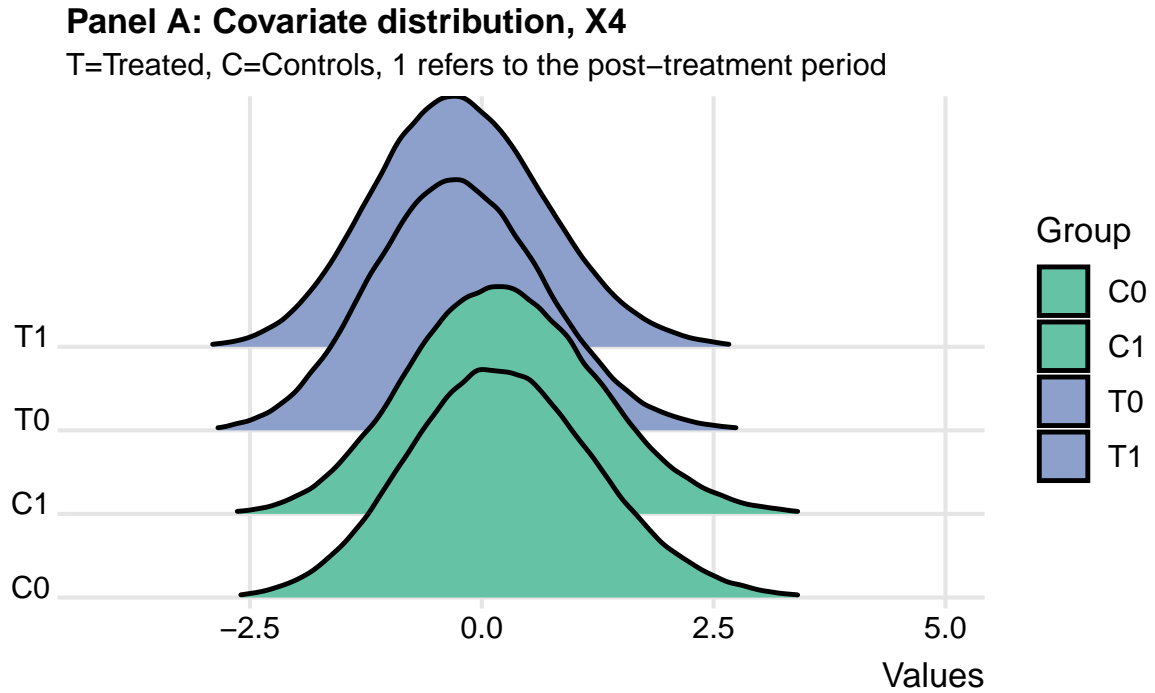
Joel Slemrod and Shlomo Yitzhaki. Tax avoidance, evasion, and administration. 3: 1423–1470, 2002.

Bret Zeldow and Laura A Hatfield. Confounding and regression adjustment in difference-in-differences. *arXiv preprint arXiv:1911.12185*, 2019.

## 6 Main Figures

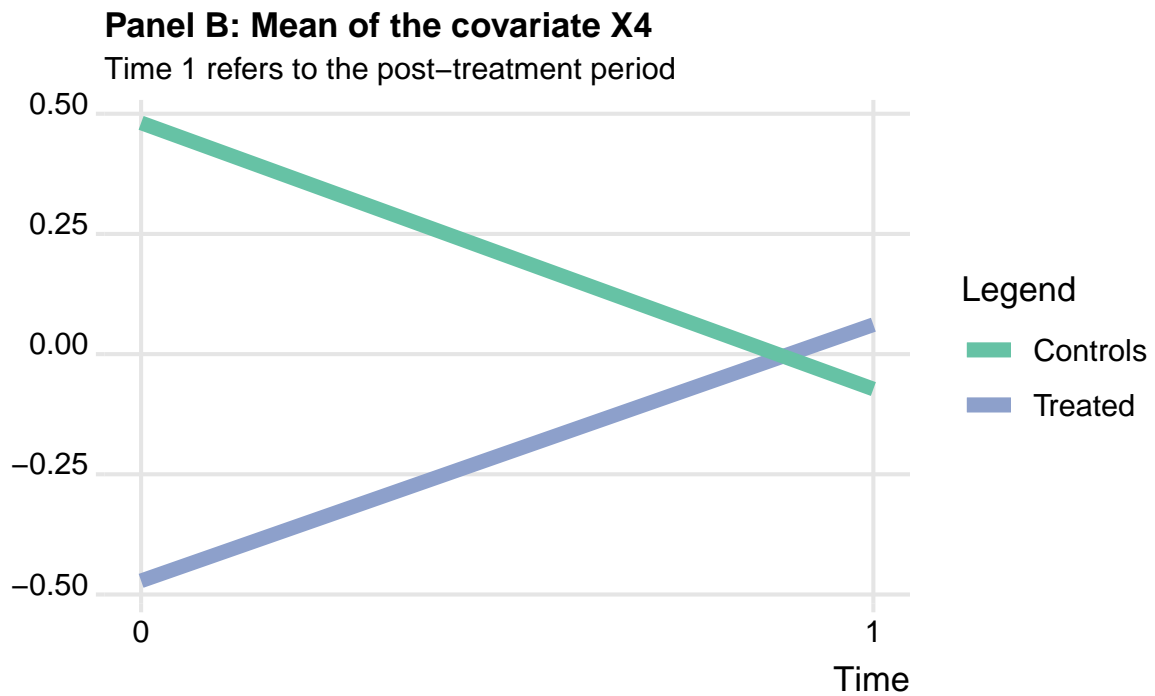
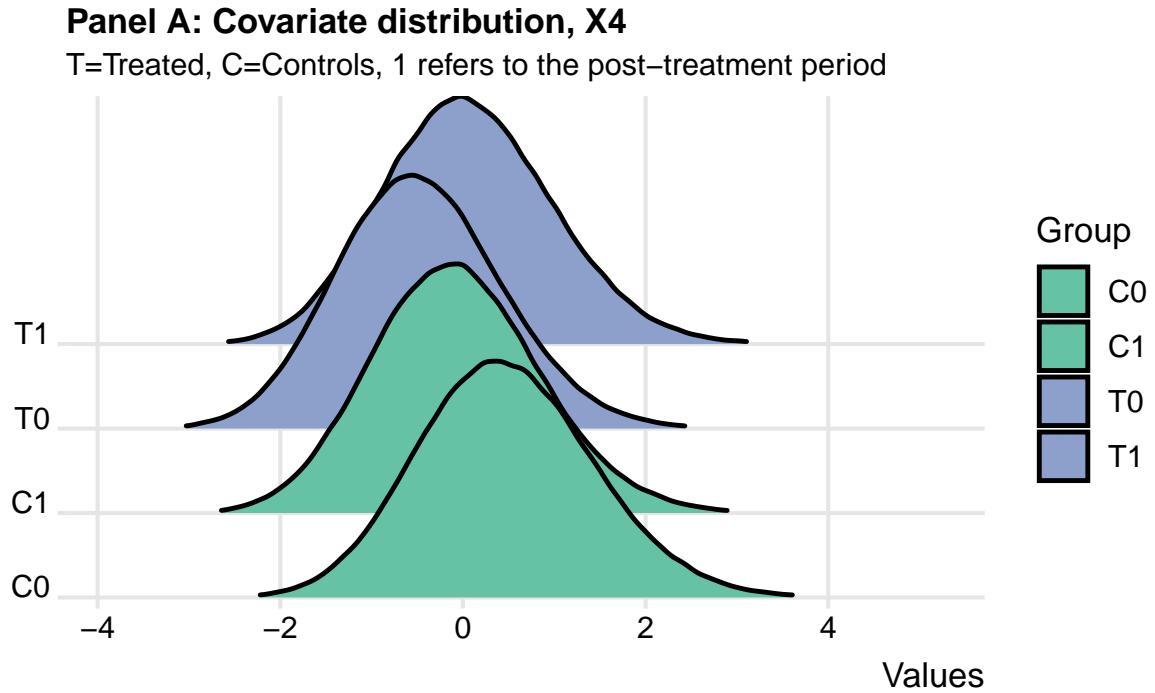


Figure 1: Distribution of  $X_4$  in Experiment 1



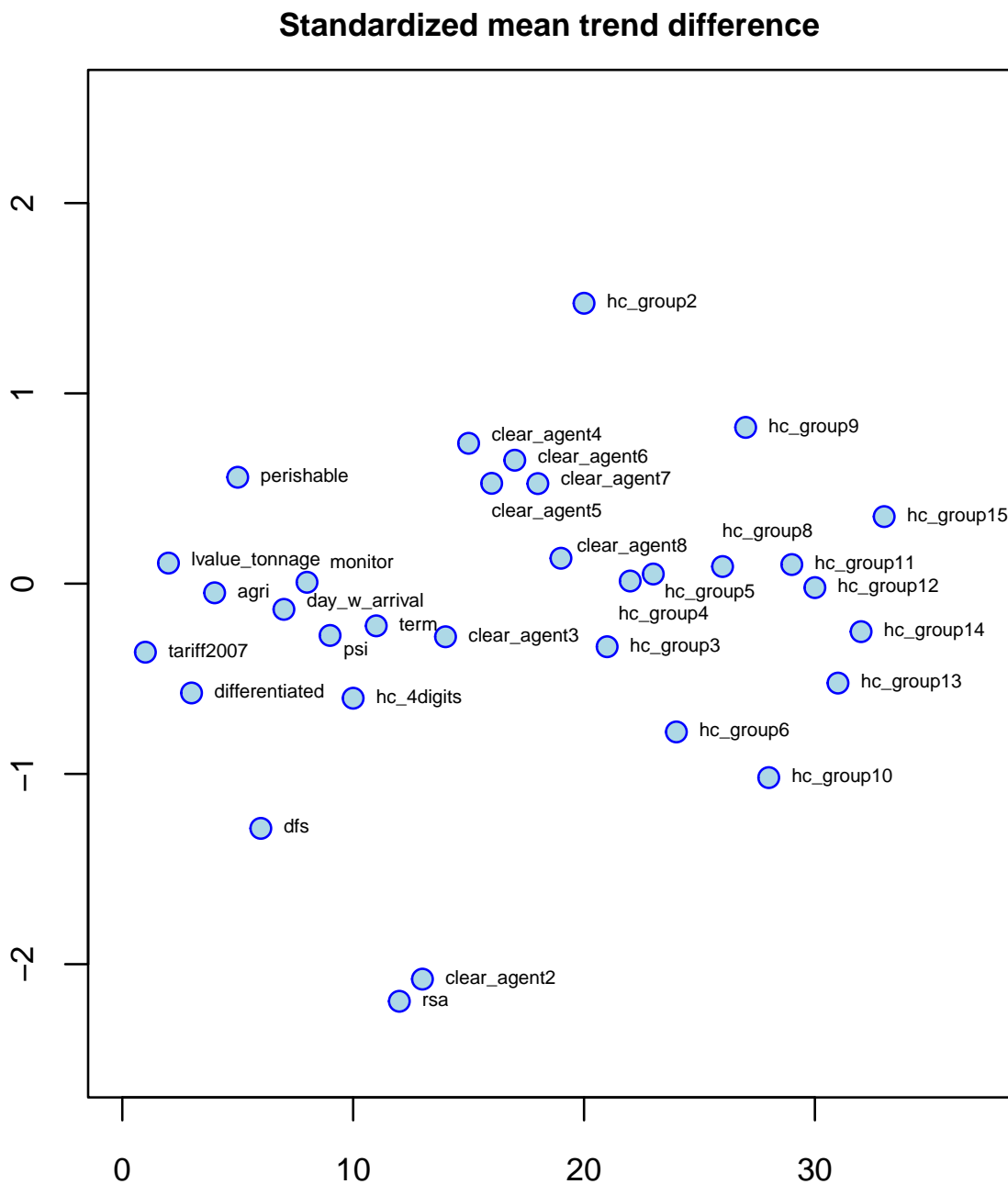
Notes: The graph considers a representative random sample from Exp.1 with DGP D. The upper plot compares the distribution of covariate  $X_4$  among the treated and controls in the pre- and post-treatment periods. The lower plot, instead, compares the respective means among treated and controls in the two time periods. Note that the distribution of  $X_4$  is time-invariant but there is heterogeneity between treated and controls populations, as captured by their difference in means.

Figure 2: Distribution of  $X_4$  in Experiment 2



Notes: The graph considers a representative random sample from Exp.2 with DGP D. The upper plot compares the distribution of covariate  $X_4$  among the treated and controls in the pre- and post-treatment periods. The lower plot, instead, compares the respective means among treated and controls in the two time periods. Note the heterogeneity in this case is also in the trend of the covariate between treated and controls.

Figure 3: Standardized mean trend difference among treated and controls for each covariate



Notes: The graph plots the 33 covariates used as controls. The covariates are enumerated and ordered arbitrarily from left to right for graphical purposes. The y-axis display the standardized mean difference in trend among treated and controls, namely  $(\bar{X}_{11} - \bar{X}_{10}) - (\bar{X}_{01} - \bar{X}_{00})$  divided by the standard deviation of X, where the overbar indicates the mean. In case of time-invariant control this measure should be 0. The figure therefore indicates the presence of a strong heterogeneity in the evolution of the covariates among treated and controls.

## 7 Main Tables

Table 1: Summary table of the estimators analyzed in the Monte Carlo simulations

Estimator	Description
TWFE	Two-Way-Fixed-Effects regression with covariates as in Eq. (1)
TWFE CORR	Two-Way-Fixed-Effects correction as in Eq. (12)
IPW	Inverse probability weighting (Abadie, 2005)
DMLDiD	Debiased machine learning IPW using lasso first-stage estimates (Chang, 2020)
DIPW	Double inverse probability weighting
OR	Outcome regression (Heckman et al., 1997)
DRDiD	Locally efficient doubly robust estimator, original version (Sant'Anna and Zhao, 2020)
DR-DIPW	Locally efficient doubly robust estimator with DIPW weights
IMP DRDiD	Improved locally efficient doubly robust estimator, original version (Sant'Anna and Zhao, 2020)
IMP DR-DIPW	Improved locally efficient doubly robust estimator with DIPW weights
LASSO DR-DIPW	Locally efficient doubly robust estimator, with DIPW weights and lasso first-stage estimates
RF DR-DIPW	Locally efficient doubly robust estimator, DIPW weights and random forest first-stage estimates

Table 2: DGPs in Experiment 1 (PS=propensity score, OR=outcome regression)

<b>DGP.A (PS and OR models correct)</b>	<b>DGP.B (PS model incorrect, OR correct)</b>
$Y_{d,0} = f_{reg}(Z) + v(Z, D) + \epsilon_0(d)$	$Y_{d,0} = f_{reg}(Z) + v(Z, D) + \epsilon_0(d)$
$Y_{d,1} = 2 \cdot f_{reg}(Z) + v(Z, D) + \epsilon_1(d)$	$Y_{d,1} = 2 \cdot f_{reg}(Z) + v(Z, D) + \epsilon_1(d)$
$p(Z) = \frac{\exp(f_{ps}(Z))}{(1 + \exp(f_{ps}(Z)))}$	$p(X) = \frac{\exp(f_{ps}(X))}{(1 + \exp(f_{ps}(X)))}$
$\lambda = 0.5$	$\lambda = 0.5$
$D = 1\{p(Z) \geq U_d\}$	$D = 1\{p(X) \geq U_d\}$
$T = 1\{\lambda \geq U_t\}$	$T = 1\{\lambda \geq U_t\}$
<b>DGP.C (PS model correct, OR incorrect)</b>	<b>DGP.D (PS and OR models incorrect)</b>
$Y_{d,0} = f_{reg}(X) + v(X, D) + \epsilon_0(d)$	$Y_{d,0} = f_{reg}(X) + v(X, D) + \epsilon_0(d)$
$Y_{d,1} = 2 \cdot f_{reg}(X) + v(X, D) + \epsilon_1(d)$	$Y_{d,1} = 2 \cdot f_{reg}(X) + v(X, D) + \epsilon_1(d)$
$p(Z) = \frac{\exp(f_{ps}(Z))}{(1 + \exp(f_{ps}(Z)))}$	$p(X) = \frac{\exp(f_{ps}(X))}{(1 + \exp(f_{ps}(X)))}$
$\lambda = 0.5$	$\lambda = 0.5$
$D = 1\{p(Z) \geq U_d\}$	$D = 1\{p(X) \geq U_d\}$
$T = 1\{\lambda \geq U_t\}$	$T = 1\{\lambda \geq U_t\}$

Notes: EXP.1 assumes a non-randomized experiment, homogeneous effects in X and time-invariant covariates.

Table 3: Exp.1A Propensity score model correct, outcome regression correct

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	20.762	21.071	12.952	0.002
TWFE CORR	Regression, Eq. (3)	0.002	0.201	0.040	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	0.158	9.587	91.892	0.005
DMLDiD	<a href="#">Chang (2020)</a>	34.356	71.469	3,927.540	43.660
DIPW	Author’s work	0.002	4.705	22.141	0.009
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	0.102	7.585	57.516	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.001	0.210	0.044	0.011
DR-DIPW	Author’s work	0.001	0.211	0.044	0.015
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.001	0.210	0.044	0.015
IMP DR-DIPW	Author’s work	0.001	0.210	0.044	0.025
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	0.123	0.347	0.105	2.768
RF DR-DIPW	Author’s work	1.618	3.450	9.286	6.097

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.1 assumes a non-randomized experiment, homogeneous effects in X, and time-invariant covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between D, T, and X (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.

Table 4: Exp.1B Propensity score model incorrect, outcome regression model correct

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	19.139	19.490	13.569	0.002
TWFE CORR	Regression, Eq. (3)	0.004	0.203	0.041	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	0.854	9.793	95.167	0.005
DMLDiD	<a href="#">Chang (2020)</a>	81.060	109.279	5,371.175	43.882
DIPW	Author’s work	0.833	3.980	15.146	0.009
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	0.031	8.196	67.179	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.005	0.210	0.044	0.011
DR-DIPW	Author’s work	0.005	0.210	0.044	0.016
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.005	0.211	0.045	0.015
IMP DR-DIPW	Author’s work	0.005	0.211	0.045	0.025
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	0.007	0.341	0.116	3.003
RF DR-DIPW	Author’s work	1.316	3.621	11.380	6.395

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.1 assumes a non-randomized experiment, homogeneous effects in X, and time-invariant covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between D, T, and X (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.

Table 5: Exp.1C Propensity score model correct, outcome regression model incorrect

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	13.117	14.056	25.524	0.002
TWFE CORR	Regression, Eq. (3)	1.291	4.715	20.565	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	0.030	9.204	84.717	0.005
DMLDiD	<a href="#">Chang (2020)</a>	62.061	99.318	6,012.513	44.992
DIPW	Author’s work	0.054	5.575	31.082	0.009
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	1.444	8.079	63.188	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.004	4.748	22.543	0.011
DR-DIPW	Author’s work	0.004	4.665	21.761	0.016
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.061	4.079	16.637	0.015
IMP DR-DIPW	Author’s work	0.068	4.086	16.690	0.025
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	0.143	3.454	11.908	3.029
RF DR-DIPW	Author’s work	0.055	3.607	13.009	6.465

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.1 assumes a non-randomized experiment, homogeneous effects in X, and time-invariant covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between D, T, and X (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.



Table 6: Exp.1D Propensity score model incorrect, outcome regression model incorrect

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	16.269	17.067	26.619	0.002
TWFE CORR	Regression, Eq. (3)	5.108	6.919	21.778	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	4.037	10.569	95.401	0.005
DMLDiD	<a href="#">Chang (2020)</a>	115.748	145.629	7,810.146	43.449
DIPW	Author’s work	3.915	6.940	32.833	0.009
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	5.248	10.001	72.473	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	3.166	6.048	26.558	0.011
DR-DIPW	Author’s work	3.164	5.960	25.508	0.016
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	2.550	4.874	17.258	0.015
IMP DR-DIPW	Author’s work	2.563	4.886	17.309	0.025
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	1.938	4.212	13.988	3.105
RF DR-DIPW	Author’s work	2.699	4.770	15.470	6.430

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.1 assumes a non-randomized experiment, homogeneous effects in X, and time-invariant covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between D, T, and X (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.

Table 7: DPGs in Experiment 2 (PS=propensity score, OR=outcome regression)

<b>DGP.A (PS and OR models correct)</b>	<b>DGP.B (PS model incorrect, OR correct)</b>
$Y_{d,0} = f_{reg}(Z) + v(Z, D) + \epsilon_0(d)$	$Y_{d,0} = f_{reg}(Z) + v(Z, D) + \epsilon_0(d)$
$Y_{d,1} = 2 \cdot f_{reg}(Z) + v(Z, D) + \delta(Z) \cdot D + \epsilon_1(D)$	$Y_{d,1} = 2 \cdot f_{reg}(Z) + v(Z, D) + \delta(Z) \cdot D + \epsilon_1(D)$
$p(Z) = \frac{\exp(f_{ps}(Z))}{(1 + \exp(f_{ps}(Z)))}$	$p(X) = \frac{\exp(f_{ps}(X))}{(1 + \exp(f_{ps}(X)))}$
$t(D, Z) = D \cdot p(-Z) + (1 - D) \cdot p(Z)$	$t(D, X) = D \cdot p(-X) + (1 - D) \cdot p(X)$
$D = 1\{p(Z) \geq U_d\}$	$D = 1\{p(X) \geq U_d\}$
$T = 1\{\lambda(Z) \geq U_t\}$	$T = 1\{\lambda(X) \geq U_t\}$
<b>DGP.C (PS model correct, OR incorrect)</b>	<b>DGP.D (PS and OR models incorrect)</b>
$Y_{d,0} = f_{reg}(X) + v(X, D) + \epsilon_0(d)$	$Y_{d,0} = f_{reg}(X) + v(X, D) + \epsilon_0(d)$
$Y_{d,1} = 2 \cdot f_{reg}(X) + v(X, D) + \delta(X) \cdot D + \epsilon_1(D)$	$Y_{d,1} = 2 \cdot f_{reg}(X) + v(X, D) + \delta(X) \cdot D + \epsilon_1(D)$
$p(Z) = \frac{\exp(f_{ps}(Z))}{(1 + \exp(f_{ps}(Z)))}$	$p(X) = \frac{\exp(f_{ps}(X))}{(1 + \exp(f_{ps}(X)))}$
$t(D, Z) = D \cdot p(-Z) + (1 - D) \cdot p(Z)$	$t(D, X) = D \cdot p(-X) + (1 - D) \cdot p(X)$
$D = 1\{p(Z) \geq U_d\}$	$D = 1\{p(X) \geq U_d\}$
$T = 1\{\lambda(Z) \geq U_t\}$	$T = 1\{\lambda(X) \geq U_t\}$

Notes: EXP.2 assumes a non-randomized experiment, heterogeneous effects in X and time-varying covariates.

Table 8: 2A Propensity score model correct, outcome regression model correct

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	8.928	9.695	14.276	0.002
TWFE CORR	Regression, Eq. (3)	0.002	0.219	0.048	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	45.102	46.124	93.208	0.005
DMLDiD	<a href="#">Chang (2020)</a>	297.085	316.542	11, 939.030	43.685
DIPW	Author’s work	0.481	6.689	44.510	0.010
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	26.066	27.144	57.363	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.002	0.226	0.051	0.011
DR-DIPW	Author’s work	0.001	0.266	0.071	0.016
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.001	0.237	0.056	0.015
IMP DR-DIPW	Author’s work	0.001	0.259	0.067	0.026
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	0.451	0.538	0.086	3.379
RF DR-DIPW	Author’s work	6.172	7.191	13.619	6.511

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.2 assumes a non-randomized experiment, heterogeneous effects in  $X$  and time-varying covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between  $D$ ,  $T$ , and  $X$  (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.

Table 9: 2B Propensity score model incorrect, outcome regression model correct

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	9.165	9.875	13.518	0.002
TWFE CORR	Regression, Eq. (3)	0.003	0.217	0.047	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	51.096	52.055	98.854	0.005
DMLDiD	<a href="#">Chang (2020)</a>	308.232	332.434	15, 505.580	43.200
DIPW	Author’s work	4.567	19.656	365.507	0.010
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	32.182	33.184	65.478	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.003	0.221	0.049	0.011
DR-DIPW	Author’s work	0.002	0.256	0.066	0.016
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	0.003	0.232	0.054	0.015
IMP DR-DIPW	Author’s work	0.004	0.249	0.062	0.027
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	0.241	0.387	0.092	3.590
RF DR-DIPW	Author’s work	6.393	7.600	16.884	6.775

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.2 assumes a non-randomized experiment, heterogeneous effects in  $X$  and time-varying covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between  $D$ ,  $T$ , and  $X$  (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.

Table 10: 2C Propensity score model correct, outcome regression model incorrect

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	5.816	7.863	28.010	0.002
TWFE CORR	Regression, Eq. (3)	4.966	6.720	20.499	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	31.208	32.585	87.807	0.005
DMLDiD	<a href="#">Chang (2020)</a>	338.258	358.140	13,845.390	44.272
DIPW	Author’s work	0.296	6.699	44.785	0.010
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	21.740	23.115	61.659	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	4.379	6.437	22.267	0.011
DR-DIPW	Author’s work	0.277	5.743	32.910	0.016
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	1.092	4.562	19.616	0.015
IMP DR-DIPW	Author’s work	0.519	4.569	20.604	0.027
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	0.829	3.797	13.728	3.620
RF DR-DIPW	Author’s work	0.127	3.867	14.934	6.826

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.2 assumes a non-randomized experiment, heterogeneous effects in  $X$  and time-varying covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between  $D$ ,  $T$ , and  $X$  (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.

Table 11: 2D Propensity score model incorrect, outcome regression model incorrect

Estimator	Reference	Bias	RMSE	Variance	Time
<b>TWFE</b>					
TWFE	Regression, Eq. (1)	3.437	6.240	27.123	0.002
TWFE CORR	Regression, Eq. (3)	16.414	17.066	21.834	0.002
<b>IPW</b>					
IPW	<a href="#">Abadie (2005)</a>	54.838	55.659	90.716	0.005
DMLDiD	<a href="#">Chang (2020)</a>	342.449	367.046	17,451.260	43.129
DIPW	Author’s work	7.553	14.246	145.889	0.009
<b>OR</b>					
OR	<a href="#">Heckman et al. (1997)</a>	44.137	44.880	66.163	0.004
<b>Doubly-Robust</b>					
DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	17.736	18.416	24.575	0.011
DR-DIPW	Author’s work	13.740	16.231	74.672	0.016
IMP DRDiD	<a href="#">Sant’Anna and Zhao (2020)</a>	12.793	13.643	22.459	0.015
IMP DR-DIPW	Author’s work	10.429	11.551	24.664	0.026
<b>Debiased ML</b>					
LASSO DR-DIPW	Author’s work	7.586	8.898	21.640	3.674
RF DR-DIPW	Author’s work	11.321	12.142	19.274	6.780

Notes: Simulations based on sample size  $n = 1000$  and 10000 Monte Carlo repetitions. EXP.2 assumes a non-randomized experiment, heterogeneous effects in  $X$  and time-varying covariates. TWFE is the standard regression specification with naively adding a set of covariates (Eq. (1)); TWFE CORR is the regression correction that adds also all possible interaction terms between  $D$ ,  $T$ , and  $X$  (Eq. (3)); IPW is the inverse probability weighting (Eq. (6)); DMLDiD is the debiased machine learning version of the IPW estimator using lasso; DIPW is the double inverse probability weighting estimator (Eq. (10)); DRDiD is the locally-efficient doubly robust estimator as in (Eq. (9)) and it is proposed in its “improved” version IMP DRDiD; likewise, DR-DIPW is the locally-efficient doubly robust estimator with DIPW weights (Eq. (12)), which is also proposed in its “improved” (IMP DR-DIPW), lasso (LASSO DR-DIPW) and random forest (RF DR-DIPW) versions. If not otherwise specified, the propensity score is estimated with logit and the outcome model through linear regression. Finally, ‘Bias’, ‘RMSE’, ‘Variance’, and ‘Time’, stand for the average simulated absolute bias, simulated root mean-squared errors, average estimator variance, and average required computational time respectively. Refer to the main text for further details.

Table 12: Variables included in  $\Gamma_i$ 

	Description
diff	If the product have differentiated prices among countries
agri	If the product is an agricultural good
lvalue	The log shipment value per tonnage
perishable	If the product is perishable
largefirm	If the firm has has more than 100 employees
dayarrival	The day of arrival during the week
inspection	If the shipment was pre-inspected at origin
monitor	If the shipment was monitored
SouthAfrica	If the product comes from South Africa
terminal	Terminal of clearence
hs4group	4-digits Harmonized System (HS) code for product industry classification

Table 13: The effect of tariff reduction on bribes

	TWFE <a href="#">Sequeira (2016)</a>	TWFE ( $\Gamma_i \cdot POST$ ) <a href="#">Sequeira (2016)</a>	DMLDiD (Kernel) <a href="#">Chang (2020)</a>	DMLDiD (lasso) <a href="#">Chang (2020)</a>
ATT	-3.748***	-2.928***	-6.998*	-5.222**
St.Err.	1.075	0.944	3.752	2.647

Notes: TWFE and TWFE( $\Gamma_i \times POST$ ) are Equation 1 and 2 in [Sequeira \(2016\)](#): the first controls for covariates, while the second adds also the interactions between covariates and the post-treatment dummy. DMLDiD (Kernel) and DMLDiD (lasso) are Column 3 and 5 in Table 2 in [Chang \(2020\)](#). Since the estimator is an IPW method adapted to handle machine-learning first stage estimates, the first uses Kernel in the first stage while the latter employs lasso. The coefficients capture the difference in the log of bribes paid for products that changed tariff level, before and after the tariff change took place. Standard errors are clustered at the level of product's four-digit HS code. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

Table 14: The effect of tariff reduction on bribes

	LASSO DR-DIPW Assumption 5.a	LASSO DR-DIPW Assumption 5.b
Coefficient	-2.478**	-1.594
St.Err.	1.052	1.052

Notes: LASSO DR-DIPW is Eq. (12) with lasso first stage estimates, where the standard errors are computed through a bootstrap procedure. Assumption 5.a imposes the absence of bad controls, while Assumption 5.b asserts that treated and controls would have had a parallel trend in case the treated had not received the treatment. The coefficients capture the difference in the log of bribes paid for products that changed tariff level, before and after the tariff change took place. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .