

Learning the effect of persuasion via difference-in-differences

Sung Jae Jun
Sokbae Lee

The Institute for Fiscal Studies
Department of Economics, UCL

cemmap working paper CWP24/24



LEARNING THE EFFECT OF PERSUASION VIA DIFFERENCE-IN-DIFFERENCES*

SUNG JAE JUN[†]

SOKBAE LEE[‡]

Pennsylvania State University Columbia University

December 5, 2024

Abstract. The persuasion rate is a key parameter for measuring the causal effect of a directional message on influencing the recipient's behavior. Its identification has relied on exogenous treatment or the availability of credible instruments, but the requirements are not always satisfied in observational studies. Therefore, we develop a novel econometric framework for the average persuasion rate on the treated and other related parameters by using the difference-in-differences approach. The average treatment effect on the treated is a standard parameter in difference-in-differences, but we show that it is an overly conservative measure in the context of persuasion. For estimation and inference, we propose regression-based approaches as well as semiparametrically efficient estimators. Beginning with the two-period case, we extend the framework to staggered treatment settings, where we show how to conduct richer analyses like the event-study design. We investigate the British election and the Chinese curriculum reform as empirical examples.

Key Words: Persuasion Rate, Treatment Effects, Media, Voting Behavior, Event Study Design.

JEL Classification Codes: C22, C23, D72, L82

*We would like to thank Peng Ding, Stefano DellaVigna, Vitor Possebom, Jonathan Roth, Pedro Sant'Anna, Zeyang (Arthur) Yu, and the seminar participants at McMaster University, Munich Econometrics Seminar, Seoul National University, Sungkyunkwan University, the University of Kentucky, and the University of Toronto for encouraging and helpful comments.

[†]Department of Economics, suj14@psu.edu.

[‡]Department of Economics, sl3841@columbia.edu.

1. INTRODUCTION

Since the dawn of civilization, the meaning of persuasion has evolved from Aristotle’s *Rhetoric* to Jane Austen’s *Persuasion* and to curriculum reform in China (e.g., [Cantoni et al., 2017](#)). Persuasion is integral to both democracy and the market economy. Economists have empirically measured how persuasive efforts influence the actions of consumers, voters, donors, and investors (see [DellaVigna and Gentzkow, 2010](#), for a survey of the early literature).

One of the popular measures for the effectiveness of persuasive efforts is the persuasion rate that was first used by [DellaVigna and Kaplan \(2007\)](#) in the context of a biased news media and its influence on voters’ behaviors. It has been regarded as a standardized measure to compare the magnitudes of treatment effects across different settings. For example, Table 1 in [DellaVigna and Gentzkow \(2010\)](#) and Figure 7 in [Bursztyn and Yang \(2022\)](#)¹ show the estimates of persuasion rates across many papers. [Jun and Lee \(2023\)](#) then reformulate the persuasion rate as a causal parameter to measure the effect of a persuasive message on influencing the recipient’s behavior. They conduct a formal identification analysis within the potential outcome framework, for which they largely rely on exogenous treatment or the availability of credible instruments. However, those requirements are not always satisfied in observational studies. Therefore, we propose the average persuasion rate on the treated (APRT) and its reverse version (R-APRT) as causal parameters of interest, and we develop a framework to identify, estimate, and conduct inference for them via difference-in-differences (DID). The average treatment effect on the treated (ATT) is a popular causal parameter in this context, but we show that it is an overly conservative measure compared to APRT and R-APRT.

To convey the idea, consider the following example. Let D be a binary indicator for an exposure to a certain media platform that delivers a directional message favoring a particular political party. For $d \in \{0, 1\}$, let $Y(d)$ be a binary potential outcome that indicates

¹Specifically, [Bursztyn and Yang \(2022\)](#) provide a meta-analysis of the literature that shows that experimental treatments to recalibrate misperceptions can lead to changes in behaviors. Figure 7 in [Bursztyn and Yang \(2022\)](#) displays the persuasion rates based on [DellaVigna and Kaplan \(2007\)](#).

whether the agent votes for the favored party or not. Since the media platform delivers a directional message, suppose that $Y(1) \geq Y(0)$ with probability one: we refer to this assumption as the assumption of no backlash. Then, ATT can be written as

$$\mathbb{P}\{Y(1) = 1 \mid D = 1\} - \mathbb{P}\{Y(0) = 1 \mid D = 1\} = \mathbb{P}\{Y(1) = 1, Y(0) = 0 \mid D = 1\}.$$
²

However, ATT is a different measure from the proportion of those who voted for the favored party after the media exposure among the people in the treatment group who would have voted differently without the media exposure: i.e.,

$$\mathbb{P}\{Y(1) = 1, Y(0) = 0 \mid D = 1\} \leq \mathbb{P}\{Y(1) = 1 \mid Y(0) = 0, D = 1\}, \quad (1)$$

where the inequality can be severely strict when $\mathbb{P}\{Y(0) = 0 \mid D = 1\}$ is small. If the size of the target audience to persuade among the treated is small, then ATT will also be small even if the message is quite successful to change their behaviors.

The conditional probability on the right-hand side of (1) is what we call APRT, which focuses only on the subpopulation that is relevant to measure the persuasive effect of a directional message. However, this target audience is a counterfactual group that is not directly observable. Therefore, as a supplementary parameter, we also consider the reverse version of APRT (R-APRT), i.e., $\mathbb{P}\{Y(0) = 0 \mid Y(1) = 1, D = 1\}$, which shows what proportion of those who were exposed to a certain media platform, and who actually voted for the party the media endorsed would have voted differently if it were not for the media exposure. In other words, it measures the degree of necessity of the media exposure for the voting behavior desired by the media. Indeed, R-APRT is exactly what [Pearl \(1999\)](#) refers to as the probability of necessity (PN).³ However, we remark that the forward version of APRT in equation (1) is distinct from [Pearl \(1999\)](#)'s probability of sufficiency (PS), which would correspond to the (average) persuasion rate on the untreated, i.e., $\mathbb{P}\{Y(1) = 0 \mid$

²The equality follows from the fact that $\mathbb{P}\{Y(1) = 1 \mid D = 1\} = \mathbb{P}\{Y(1) = 1, Y(0) = 1 \mid D = 1\} + \mathbb{P}\{Y(1) = 1, Y(0) = 0 \mid D = 1\}$ and $\mathbb{P}\{Y(0) = 1 \mid D = 1\} = \mathbb{P}\{Y(0) = 1, Y(1) = 1 \mid D = 1\}$ under no backlash.

³The same parameter is called the probability of causal attribution in the political science literature ([Yamamoto, 2012](#)).

$Y(0) = 1, D = 0$ in our context. We do not consider PS because it will not be point-identified via DID, just like the average treatment effect on the untreated is not identifiable via DID. See section 2 for more discussion.

The idea of “rescaling” to focus on a particular event or subpopulation has been around for a long time: e.g., [Bayes \(1763\)](#) rescales joint probability to obtain conditional probability; [Imbens and Angrist \(1994\)](#) rescale the effect of the intent-to-treat (ITT) to obtain the local average treatment effect (LATE); [Jun and Lee \(2023\)](#) use rescaling the average treatment effect (ATE) and LATE to obtain the persuasion rate and the local persuasion rate, respectively. Here, we follow the same idea, and we rescale ATT to obtain APRT and R-APRT in the absence of backlash. Furthermore, even if the assumption of no backlash is violated, scaling up ATT leads to lower bounds on APRT and R-APRT. As it is the case for joint versus conditional probabilities, ATT, APRT, and R-APRT are all related, but they deliver different information in general.

We not only articulate identification issues for APRT and R-APRT in panel-like setups with no instruments, but we also investigate estimation and inference problems in depth. Specifically, the main contributions of the paper are threefold. First, we clarify the relationships between APRT, R-APRT, and ATT with and without backlash: for instance, in the absence of backlash, identification of APRT as well as that of R-APRT is reduced to that of ATT. For the latter, we follow [Abadie \(2005\)](#) and [Callaway and Sant’Anna \(2021\)](#), and we focus on the setting where the assumption of parallel trends holds after conditioning on observed covariates. We consider the canonical two-period case as well as the popular case of staggered treatment adoption. In this way, we contribute to the recent fast-growing literature on the DID framework and the panel event-study design: DID and heterogeneous treatment effects (e.g. [De Chaisemartin and d’Haultfoeuille, 2020](#); [Xu et al., 2024](#)), estimation of causal effects with panel data (e.g. [Ding and Li, 2019](#); [Freyaldenhoven et al., 2019](#); [Arkhangelsky et al., 2021](#)), and heterogeneous treatment timings (e.g. [Goodman-Bacon, 2021](#); [Sun and Abraham, 2021](#); [Callaway and Sant’Anna, 2021](#); [Athey and Imbens, 2022](#)).

See e.g., [de Chaisemartin and D’Haultfœuille \(2022\)](#); [Freyaldenhoven et al. \(2021\)](#); [Roth et al. \(2023\)](#); [Sun and Shapiro \(2022\)](#) for the introductory articles and the latest surveys.

Second, we propose two alternative estimation methods. When covariates are not needed, we show that one can use simple regression-based approaches for estimation, and we explain how they are connected to the standard two-way fixed effect regression and event-study design models. When covariates are needed, “partialing them out” from the regression analysis could be an easy option, but we do not advocate it because of potential contamination bias problems in estimating treatment effects as [Blandhol et al. \(2022\)](#) and [Goldsmith-Pinkham et al. \(2022\)](#) pointed out. Instead, we propose a battery of asymptotically equivalent efficient estimators of APRT and R-APRT, including doubly (and locally) robust estimators based on our derivation of the efficient influence functions of APRT and R-APRT. Therefore, we contribute to the literature on efficient and doubly robust estimation of treatment effects for which ATE and ATT are typical parameters of interest (see, e.g., [Hahn, 1998](#); [Hirano et al., 2003](#); [Sant’Anna and Zhao, 2020](#); [Chernozhukov et al., 2022](#), among others).

Third, we offer two options for inference in addition to the straightforward method based on the regression-based approaches. One relies on our efficient estimators, while the other involves back-of-the-envelope inference based on ATT. The former, though efficient and asymptotically exact, necessitates access to the entire dataset, whereas the latter is potentially conservative but can be conducted by combining the confidence interval for ATT with information on a certain summary statistic, eliminating the need for the entire dataset. In other words, the latter provides a convenient way to obtain the confidence intervals for APRT and R-APRT when the confidence interval for ATT is already provided.

We are not aware of any existing work that studies the persuasion rate within the DID framework. However, there are related papers on the persuasion rate in general. In the same setting as [Jun and Lee \(2023\)](#), [Yu \(2023\)](#) explores identification of the statistical characteristics of persuasion types. [Possebom and Riva \(2024\)](#) study identification of the average persuasion rate with sample selection. The concept of the causal persuasion rate is

related to probabilities of causation, which have been studied for decades in the literature outside economics: e.g., [Pearl \(1999\)](#); [Yamamoto \(2012\)](#); [Dawid et al. \(2014\)](#); [Dawid and Musio \(2022\)](#); [Zhang et al. \(2024\)](#). However, APRT does not appear to match any of the probabilities of causation, and furthermore, the existing work on them is mostly limited to an analysis of population bounds without relating to DID or investigating inferential problems. Furthermore, economists were not cognizant of this literature; one recent exception is [Possebom and Riva \(2024\)](#). More broadly speaking, our paper is related to a branch of the literature on treatment effects that depend on the joint distribution of potential outcomes. For example, see [Heckman et al. \(1997\)](#) as early pioneering work as well as [Ji et al. \(2023\)](#) and [Kaji and Cao \(2023\)](#) for recent working papers.

The remainder of the paper is organized as follows. Sections 2 to 7 focus on the canonical case of two periods, whereas section 8 shows extensions to the case of staggered treatments in detail. In section 7, we present an empirical example of the two-period case by re-examining the impact of news media on the 1997 British election ([Ladd and Lenz, 2009](#); [Hainmueller, 2012](#)). In exploring staggered treatments in section 8, we emphasize how to do an event-study-design-like analysis. In section 9, as an empirical example of staggered treatments, we re-evaluate the effects of curriculum reforms on political attitudes in China ([Cantoni et al., 2017](#)). Finally, the appendices, including the online ones, include all the proofs as well as additional results and discussions.

2. THE PARAMETERS

Let Y_t and D_t be binary random variables that are observed at time t , where $t \in \{0, 1\}$: $t = 0$ represents a pre-treatment period. We will denote the potential outcomes at time t by $Y_t(d)$, and therefore we observe $Y_t = D_t Y_t(1) + (1 - D_t) Y_t(0)$. We also observe a vector of exogenous covariates, which will be denoted by X . It will be helpful to have an example in mind: e.g., $Y_1(1)$ indicates whether an agent votes for a certain party or not, after getting exposed to a certain news platform in period 1, whereas $Y_1(0)$ is whether she votes for the party without having an exposure to the platform in period 1.

Our goal is to study identification and estimation of the causal persuasion rate on the treated in a setup of panel data: we will formally define and discuss the parameter of interest in the following subsections. We first make the following assumption to describe the setup formally. Let \mathcal{X} be the support of X .

Assumption A (No Anticipation and Non-Degenerate Probabilities). *At time $t = 0$, no one is treated so that $Y_0 = Y_0(0)$ and $D_0 = 0$ with probability one. At time $t = 1$, there is a constant $\epsilon > 0$ such that for (almost) all $x \in \mathcal{X}$, $\epsilon \leq \min[\mathbb{P}\{Y_1(0) = 0, D_1 = 1 \mid X = x\}, \mathbb{P}\{Y_1(1) = 1, D_1 = 1 \mid X = x\}]$ and $\mathbb{P}(D_1 = 1 \mid X = x) \leq 1 - \epsilon$.*

Assumption A clarifies that $t = 0$ is a pre-treatment period, and that the treatment occurs at $t = 1$. The second part of assumption A is to ensure that, for instance, there is a non-trivial group of people who are exposed to a certain news platform and who would not have voted for the party that the news platform supports if they had not been exposed to it. Otherwise, there would be nobody to “persuade” among the treated.

In the following subsections, we will define two versions of the individualized persuasion rate on the treated, and we will discuss their aggregation.

2.1. The Conditional Persuasion Rate on the Treated. We define the conditional persuasion rate on the treated (CPRT) and its reverse version (R-CPRT) in the post-treatment period as follows: for $x \in \mathcal{X}$,

$$\begin{aligned}\theta_c(x) &:= \mathbb{P}\{Y_1(1) = 1 \mid Y_1(0) = 0, D_1 = 1, X = x\}, \\ \theta_c^{(r)}(x) &:= \mathbb{P}\{Y_1(0) = 0 \mid Y_1(1) = 1, D_1 = 1, X = x\},\end{aligned}$$

which are all well-defined under assumption A.

Although the two parameters are closely related, they have different merits. Recall the voting example we mentioned earlier, where D_1 indicates an exposure to a certain media platform that endorses a particular political party. Both $\theta_c(x)$ and $\theta_c^{(r)}(x)$ focus on those who are characterized by the exogenous covariates and who are exposed to the media platform. Within this subpopulation, $\theta_c(x)$ measures the proportion of those who have

voted for the party the news media endorsed relative to those who would not have done so without the media exposure. Therefore, it measures the direct persuasive effect of the media on the exposed in that it shows how many of those who have received the persuasive message have actually changed their behavior as the message wanted them to. This is the on-the-treated version of the persuasion rate of [DellaVigna and Kaplan \(2007\)](#) and [Jun and Lee \(2023\)](#).

Although $\theta_c(x)$ provides a straightforward interpretation, it has a drawback: i.e., it conditions on a counterfactual group that is not directly observed by the researcher. The reverse rate $\theta_c^{(r)}(x)$ addresses this issue. We now consider those who have been exposed to the media and voted for the party the media endorse. We then ask how many of them would have behaved differently if it were not for the media. This type of approach is common in the context of lawsuits and legal debates.

In fact, $\theta_c^{(r)}(x)$ is exactly [Pearl \(1999\)](#)'s probability of necessity conditional on $X = x$, which was proposed to measure the degree of necessity of a causal treatment for a "successful" outcome. In the context of the voting example, it shows how necessary the media influence is for the desired voting behavior. For additional examples that frequently arise in legal contexts, see [Pearl \(1999\)](#) and [Dawid et al. \(2014\)](#), and the accompanying comments and author responses). In contrast, $\theta_c(x)$ is distinct from a conditional version of Pearl's probability of sufficiency, i.e., $\text{PS}(x) := \mathbb{P}\{Y_1(1) = 1 \mid Y_1(0) = 0, D_1 = 0, X = x\}$, which we may call a conditional persuasion rate on the untreated. Unlike $\theta_c(x)$, $\text{PS}(x)$ is conditioned on a subpopulation that is directly identifiable. However, $\text{PS}(x)$ is not identifiable in the standard natural experiment setup with a parallel trend assumption.

To further understand CPRT and R-CPRT, let $p_{st}(s) := \mathbb{P}\{Y_1(0) = s, Y_1(1) = t \mid D_1 = 1, X = x\}$ for $(s, t) \in \{0, 1\}^2$, and we can write

$$\theta_c(x) = \frac{p_{01}(x)}{p_{01}(x) + p_{00}(x)} \quad \text{and} \quad \theta_c^{(r)}(x) = \frac{p_{01}(x)}{p_{01}(x) + p_{11}(x)}. \quad (2)$$

Here, $p_{00}(x)$ regards the "never-persuadable (NP) among the treated with $X = x$ " and $p_{11}(x)$ does the "already-persuaded (AP) among the treated with $X = x$ " while $p_{01}(x)$

considers the “treatment-persuadable (TP) among the treated with $X = x$.” Therefore, CPRT $\theta_c(x)$ is the relative share of TP in the subpopulation of NP and TP among the treated with $X = x$, whereas R-CPRT $\theta_c^{(r)}(x)$ is the relative size of TP in the subgroup of AP and TP among the treated with $X = x$.

Since the potential outcomes $Y_1(0)$ and $Y_1(1)$ cannot be observed simultaneously, we cannot generally point-identify CPRT or R-CPRT without making additional assumptions. Consider the following parameters:

$$\theta_{cL}(x) := \frac{p_{01}(x) - p_{10}(x)}{p_{01}(x) + p_{00}(x)} \quad \text{and} \quad \theta_{cL}^{(r)}(x) := \frac{p_{01}(x) - p_{10}(x)}{p_{01}(x) + p_{11}(x)}. \quad (3)$$

Unlike $\theta_c(x)$ and $\theta_c^{(r)}(x)$, $\theta_{cL}(x)$ and $\theta_{cL}^{(r)}(x)$ depend only on the marginal probabilities of the potential outcomes given $D_1 = 1$ and $X = x$ because $p_{01}(x) - p_{10}(x)$ is equal to

$$\text{CATT}(x) := \mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\} - \mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\}.$$

Since $\theta_{cL}(x) \leq \theta_c(x)$ and $\theta_{cL}^{(r)}(x) \leq \theta_c^{(r)}(x)$ in general, we can view $\theta_{cL}(x)$ and $\theta_{cL}^{(r)}(x)$ as general conservative measures of $\theta_c(x)$ and $\theta_c^{(r)}(x)$, respectively. However, they are less conservative than ATT in general. Indeed, they are Fréchet-Hoeffding lower bounds on $\theta_c(x)$ and $\theta_c^{(r)}(x)$, and they are sharp based on the marginals of the potential outcomes given $D_1 = 1$ and $X = x$. See appendix S-1 for more discussion.

Of course, if $p_{10}(x) = 0$, then $\theta_{cL}(x)$ and $\theta_{cL}^{(r)}(x)$ coincide with $\theta_c(x)$ and $\theta_c^{(r)}(x)$. Below we discuss this condition in detail.

Assumption B (No Backlash). $\mathbb{P}\{Y_1(1) \geq Y_1(0) \mid D_1 = 1, X\} = 1$ *almost surely*.

Under assumption B, we have $p_{10}(x) = 0$, and hence $p_{00}(x) + p_{01}(x) + p_{11}(x) = 1$ for (almost) all $x \in \mathcal{X}$. In other words, we rule out the presence of “contrarians” (among the treated) i.e., those who always do the opposite of what the persuasive message directs; hence, the treated (with or without conditioning on X) can be decomposed into three persuasion types: i.e., NP, AP, and TP.

In the context of the voting example, assumption **B** says that if the voters (among the treated) would vote for the party the news media support without listening to them, then exposing them to the media would not change their mind. Put differently, assumption **B** assumes that a persuasive message is directional or biased, and that there is no backlash at least among the treated. A sufficient condition for assumption **B** is the monotone treatment response (MTR) assumption (Manski, 1997), which states that $Y_1(1) \geq Y_1(0)$ almost surely in our setting. The MTR assumption with binary potential outcomes was also adopted in Pearl (1999) and Jun and Lee (2023).

Lemma 1. *If assumptions **A** and **B** holds, then $\theta_c(x) = \theta_{cL}(x)$ and $\theta_c^{(r)}(x) = \theta_{cL}^{(r)}(x)$ for all $x \in \mathcal{X}$. Even if assumption **B** can be violated, we have $\theta_c(x) \geq \theta_{cL}(x)$ and $\theta_c^{(r)}(x) \geq \theta_{cL}^{(r)}(x)$ for all $x \in \mathcal{X}$ in general.*

Lemma 1 is not an identification result yet. However, it shows that ruling out a backlash effect enables us to express the conditional probabilities $\theta_c(x)$ and $\theta_c^{(r)}(x)$ by using the marginal probabilities of the potential outcomes. Moreover, even if the backlash is of concern, $\theta_{cL}(x)$ and $\theta_{cL}^{(r)}(x)$ will continue to serve as conservative measures of $\theta_c(x)$ and $\theta_c^{(r)}(x)$, respectively. Our notation using the subscript L is to emphasize that they are robust lower bounds even if assumption **B** is violated. We provide further discussion on bounds, their sharpness, and the backlash in appendix S-1.

Therefore, under assumptions **A** and **B**, the point-identification of CPRT and R-CPRT will hinge on that of CATT. Since assumption **B** ensures that $\text{CATT}(x) \geq 0$, both CPRT and R-CPRT are guaranteed to be no smaller than CATT: i.e., $\min\{\theta_{cL}(x), \theta_{cL}^{(r)}(x)\} \geq \text{CATT}(x) \geq 0$, where the first inequality is strict if $\text{CATT}(x) > 0$ and $\mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\} > 0$: e.g., if there are people who would vote for the party a news platform endorses without even listening to them, then CPRT is strictly larger than CATT. This is because CATT does not take into account the case of “preaching to the converted” in measuring the persuasive effect, whereas CPRT and R-CPRT address this issue by conditioning.

Since $Y_1(1) = Y_1$ when $D_1 = 1$ and $Y_0(0) = Y_0$, there is only one unidentified object in $\theta_{cL}(x)$ and $\theta_{cL}^{(r)}(x)$: i.e.,

$$\tau_c(x) := \mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\}.$$

Before we proceed for identification of $\tau_c(x)$, we first discuss aggregation in the following subsection.

2.2. Aggregation. The persuasion rates on the treated as aggregated versions of $\theta_c(\cdot)$ and $\theta_c^{(r)}(\cdot)$, which we call the average persuasion rate on the treated (APRT) and reverse APRT (R-APRT), respectively, are defined by

$$\begin{aligned} \theta &:= \mathbb{E}\{\theta_c(X) \mid Y_1(0) = 0, D_1 = 1\} = \mathbb{P}\{Y_1(1) = 1 \mid Y_1(0) = 0, D_1 = 1\}, \\ \theta^{(r)} &:= \mathbb{E}\{\theta_c^{(r)}(X) \mid Y_1(1) = 1, D_1 = 1\} = \mathbb{P}\{Y_1(0) = 0 \mid Y_1(1) = 1, D_1 = 1\}. \end{aligned} \quad (4)$$

The forward version θ can be compared with the average persuasion rate (APR) that [Jun and Lee \(2023\)](#) studied, i.e., $\mathbb{P}\{Y_1(1) = 1 \mid Y_1(0) = 0\}$. The relationship between APR and APRT is similar to that of ATE and ATT.

By aggregating the numerators and denominators in equation (3), we obtain the following parameters:

$$\begin{aligned} \theta_L &:= \frac{\mathbb{E}\{\text{CATT}(X) \mid D_1 = 1\}}{\mathbb{E}\{1 - \tau_c(X) \mid D_1 = 1\}} = \frac{\text{ATT}}{\text{ATT} + \mathbb{P}(Y_1 = 0 \mid D_1 = 1)}, \\ \theta_L^{(r)} &:= \frac{\mathbb{E}\{\text{CATT}(X) \mid D_1 = 1\}}{\mathbb{E}\{\mathbb{P}(Y_1 = 1 \mid D_1 = 1, X) \mid D_1 = 1\}} = \frac{\text{ATT}}{\mathbb{P}(Y_1 = 1 \mid D_1 = 1)}, \end{aligned}$$

where $\text{ATT} := \mathbb{E}\{\text{CATT}(X) \mid D_1 = 1\}$. If there is no backlash for all $x \in \mathcal{X}$, then there is no difference between $(\theta, \theta^{(r)})$ and $(\theta_L, \theta_L^{(r)})$.

Lemma 2. *If assumptions A and B hold, then we have $\theta = \theta_L$ and $\theta^{(r)} = \theta_L^{(r)}$. Even if assumption B can be violated, we have $\theta \geq \theta_L$ and $\theta^{(r)} \geq \theta_L^{(r)}$ in general.*

It is a consequence of the Bayes rule that averaging the numerators and the denominators of $\theta_{cL}(\cdot)$ and $\theta_{cL}^{(r)}(\cdot)$ separately is the right way of aggregation. Again, identification of $\tau_c(\cdot)$ is sufficient for that of θ_L and $\theta_L^{(r)}$.

As in the case of CPRT and R-CPRT, the subscript L emphasizes that θ_L and $\theta_L^{(r)}$ provide valid lower bounds on θ and $\theta^{(r)}$ even without assumption **B**. However, if assumption **B** does not hold, then $\theta_{cL}(x)$ and $\theta_{cL}^{(r)}(x)$ can be negative for some $x \in \mathcal{X}$, whereas $\theta_c(x)$ and $\theta_c^{(r)}(x)$ can never be. Therefore, the aggregated parameters θ_L and $\theta_L^{(r)}$ may not be the best bounds we can obtain from the marginals of the potential outcomes. We discuss the issues of sharp bounds and robust interpretation without assumption **B** in detail in Online Appendix **S-1**.

For the sake of intuition, suppose that $Y_1(1) \geq Y_1(0)$ with probability one so that assumption **B** is satisfied, and focus on APRT. Then, APRT is obtained by rescaling ATT, where the rescaling factor is determined by how many people in the treatment group are not pre-converted and have potential to be persuaded:⁴ the more people in the treatment group were going to take the action of interest without listening to a persuasive message, the more important it is to address the issue of “preaching to the converted” in measuring the pure persuasive effect of the treatment. For instance, if no one in the treatment group is pre-converted and everybody is a real target to persuade, then there will be no difference between ATT and APRT. Otherwise, the two causal parameters are distinct, and the latter is usually larger.

3. IDENTIFICATION

3.1. Identification via Parallel Trends. We continue to focus on the simple case of two time periods and take a difference-in-differences (DID) approach for identification of $\tau_c(\cdot)$. The case of staggered treatment will be discussed in section **8**. The key assumption for identification is that of parallel trends.

Assumption C (Parallel Trends). $\mathbb{P}\{Y_t(0) = 1 \mid D_1 = d, X = x\}$ is separable into the sum of a time component and a treatment component: i.e., there exist functions G and H such that $\mathbb{P}\{Y_t(0) = 1 \mid D_1 = d, X = x\} = G(t, x) + H(d, x)$.

⁴In terms of the decomposition of the treatment group we discussed earlier, APRT is equal to $\#TP / (\#TP + \#NP) = \#TP / \#AP^c$, where $\#A$ represents the size of group A .

By Proposition 3.2 and Example 1 of [Roth and Sant'Anna \(2023\)](#), assumption **C** is an equivalent form of parallel trends such that

$$\begin{aligned} & \mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\} - \mathbb{P}\{Y_0(0) = 1 \mid D_1 = 1, X = x\} \\ &= \mathbb{P}\{Y_1(0) = 1 \mid D_1 = 0, X = x\} - \mathbb{P}\{Y_0(0) = 1 \mid D_1 = 0, X = x\}. \end{aligned} \quad (5)$$

The class of generalized linear models such as logit or probit is popular in parametric approaches, but assumption **C** does not allow it. For example, a parametric model such as $\mathbb{P}\{Y_t(0) = 1 \mid D_1 = d, X = x\} = \Lambda^{-1}(\beta_0 + \beta_1 t + \beta_2 d + \beta_3^\top x)$ with $\Lambda^{-1}(s) := \exp(s) / \{1 + \exp(s)\}$ is not allowed. However, it is straightforward to modify assumption **C** to introduce a nonlinear link function Λ on $[0, 1]$, as long as the link function is pre-specified.⁵ We discuss this modification in Online Appendix [S-2](#).⁶

Define

$$\Psi(X) := \Pi_0(1, X) + \Pi_1(0, X) - \Pi_0(0, X),$$

where $\Pi_t(d, x) := \mathbb{P}(Y_t = 1 \mid D_1 = d, X = x)$ for $t = \{0, 1\}$ and $(d, x^\top)^\top \in \{0, 1\} \times \mathcal{X}$, which are all directly identified from the distribution of $(Y_0, Y_1, D_1, X^\top)^\top$.

Theorem 1. *Suppose that assumptions **A** and **C** hold. Then, for all $x \in \mathcal{X}$, $\tau_c(x)$ is point-identified by $\Psi(x)$.*

Therefore, under assumptions **A** to **C**, $\theta_{cL}(\cdot)$, $\theta_{cL}^{(r)}(\cdot)$, θ_L , and $\theta_L^{(r)}$ are all point-identified: e.g., $\theta_{cL}(x)$ is identified by rescaling the usual DID parameter in that for all $x \in \mathcal{X}$,

$$\theta_{cL}(x) = \frac{\Pi_1(1, x) - \Psi(x)}{1 - \Psi(x)} = \frac{\text{CATT}(x)}{\text{CATT}(x) + 1 - \Pi_1(1, x)},$$

⁵The change-in-changes (CIC) approach of [Athey and Imbens \(2006\)](#) is popular to handle nonlinearity under alternative assumptions. However, it is less convenient when the outcomes are only binary, and therefore we do not purpose this possibility in this paper.

⁶Specifically, we consider a generalized version of parallel trends with the transformation Λ :

$$\begin{aligned} & \Lambda[\mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\}] - \Lambda[\mathbb{P}\{Y_0(0) = 1 \mid D_1 = 1, X = x\}] \\ &= \Lambda[\mathbb{P}\{Y_1(0) = 1 \mid D_1 = 0, X = x\}] - \Lambda[\mathbb{P}\{Y_0(0) = 1 \mid D_1 = 0, X = x\}]. \end{aligned}$$

where $\text{CATT}(x) = \Pi_1(1, x) - \Psi(x)$ is the usual DID estimand: i.e., $\text{CATT}(x) = \Delta(1, x) - \Delta(0, x)$ with $\Delta(d, x) := \Pi_1(d, x) - \Pi_0(d, x)$.

Since the aggregated parameters are of particular interest when we have many covariates, we make a formal statement about them as a corollary. Define

$$\bar{\theta}_L := \frac{\mathbb{E}\{\Pi_1(1, X) - \Psi(X) \mid D_1 = 1\}}{\mathbb{E}\{1 - \Psi(X) \mid D_1 = 1\}}, \quad \bar{\theta}_L^{(r)} := \frac{\mathbb{E}\{\Pi_1(1, X) - \Psi(X) \mid D_1 = 1\}}{\mathbb{E}\{\Pi_1(1, X) \mid D_1 = 1\}}.$$

Corollary 1. *Suppose that assumptions **A** to **C** hold. Then, $\theta = \theta_L = \bar{\theta}_L$ and $\theta^{(r)} = \theta_L^{(r)} = \bar{\theta}_L^{(r)}$. If assumption **B** can be violated, we have $\theta \geq \theta_L = \bar{\theta}_L$ and $\theta^{(r)} \geq \theta_L^{(r)} = \bar{\theta}_L^{(r)}$ in general.*

Therefore, if there is no backlash in the sense of assumption **B**, then APRT and R-APRT are point-identified by rescaling the usual DID estimand. Even if backlash effects are concerning, rescaling the DID parameter provides conservative measures of APRT and R-APRT. However, $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$ are not generally sharp identified lower bounds. See Online Appendix **S-1** for more details.

Recall that assumption **B** decomposes the group of the treated (with or without conditioning on X) into three persuasion types. If assumption **C** holds in addition, then ATT is identified. Therefore, the shares of the three types among the treated are all identified as well: the share of NP among the treated is given by $\mathbb{E}\{p_{00}(X) \mid D_1 = 1\} = \mathbb{P}(Y_1 = 0 \mid D_1 = 1)$, and the share of AP among the treated is $\mathbb{E}\{p_{11}(X) \mid D_1 = 1\} = \mathbb{P}(Y_1 = 1 \mid D_1 = 1) - \text{ATT}$, while that of TP among the treated is just ATT.

For estimation, we can simply replace $\Pi_t(d, x)$ with their parametric or nonparametric estimates. However, when it comes to APRT and R-APRT, directly plugging in $\Pi_t(d, x)$'s and aggregating them is not the only possibility. We will discuss various approaches for estimation in section 4.

3.2. Controlling for the Pre-Treatment Outcome and Unconfoundedness. As an alternative to the parallel trend assumption, it is popular to assume unconfoundedness after controlling for enough covariates. Specifically, when a pre-treatment outcome Y_0 is observed,

it is a natural idea to assume unconfoundedness after controlling for Y_0 in addition to X to achieve identification.

In this context, it is worth noting that if we use $Z := [Y_0, X]$ in lieu of X in defining $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$, then we are led to estimands that identify θ_L and $\theta_L^{(r)}$, respectively, under the independence assumption of $Y_1(0)$ and D_1 given Z . Therefore, adding Y_0 to X and using the DID formulas can be thought of as an implementation of identifying APRT and R-APRT via unconfoundedness given Z .

There has been some debate about the desirability of conditioning on Y_0 when estimating ATT via DID: see, e.g., [Roth et al. \(2023, pp. 2232\)](#) and the references therein. However, it seems less well-known that there is a testable condition under which it becomes moot to distinguish the two approaches. In fact, it can be shown that if Y_0 is independent of D_1 given X , then the two alternative identification assumptions lead to the same estimands for θ_L and $\theta_L^{(r)}$. Therefore, under the aforementioned independence condition, it does not matter which stance the researcher takes between the two alternative identification assumptions. We provide a more detailed discussion on this issue in appendix [S-3](#).

4. ESTIMATION

4.1. Regression-Based Approaches. It is a popular practice to estimate ATT by using a two-way fixed effect regression model. There is a similar approach to estimate the persuasion rate. In order to clarify the idea, we focus on the case where there are no covariates for now. If assumptions [A](#) to [C](#) hold without the covariates X , then APRT and R-APRT are given by

$$\bar{\theta}_L = \frac{\Delta(1) - \Delta(0)}{\Delta(1) - \Delta(0) + 1 - \Pi_1(1)} = \frac{\mathbb{E}(Y_1 - Y_0 \mid D_1 = 1) - \mathbb{E}(Y_1 - Y_0 \mid D_1 = 0)}{\mathbb{E}(1 - Y_0 \mid D_1 = 1) - \mathbb{E}(Y_1 - Y_0 \mid D_1 = 0)}, \quad (6)$$

$$\bar{\theta}_L^{(r)} = \frac{\Delta(1) - \Delta(0)}{\Pi_1(1)} = \frac{\mathbb{E}(Y_1 - Y_0 \mid D_1 = 1) - \mathbb{E}(Y_1 - Y_0 \mid D_1 = 0)}{\mathbb{E}(Y_1 \mid D_1 = 1)}, \quad (7)$$

where $\Delta(d) := \Pi_1(d) - \Pi_0(d)$ and $\Pi_t(d) := \mathbb{P}(Y_t = 1 \mid D_1 = d)$ for $(t, d) = \{0, 1\}^2$. Below we treat $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$ in (6) and (7) as the estimands of interest, for which we directly assume

that $\mathbb{E}(Y_1 \mid D_1 = 1) > 0$ and $\Psi := \mathbb{E}(Y_0 \mid D_1 = 1) + \mathbb{E}(Y_1 - Y_0 \mid D_1 = 0) < 1$, and we show that they can be obtained from a simple linear regression.

Consider the following two-way fixed effect regression model:

$$Y_{it} = \gamma_0 + G_i\gamma_1 + t\gamma_2 + tG_i\gamma + \epsilon_{it}, \quad (8)$$

where G_i is a group indicator, taking a value of 1 if individual i is in the treatment group, and 0 otherwise, and γ_1 and γ_2 capture group and time fixed effects, respectively. We make the following assumption.

Assumption D (Two-Way Fixed Effects). *(i) All those (and only those) in the treatment group are treated in period 1 so that $D_{it} = t \cdot G_i$.*

(ii) The time and group assignment is exogenous in that $\mathbb{E}(\epsilon_{it} \mid G_i, t = 0) = 0$ and $\mathbb{E}(\epsilon_{it} \mid G_i, t = 1) = 0$.

Under assumption **D**, $\gamma_0, \gamma_1, \gamma_2$, and γ can be estimated by the ordinary least squares (OLS). In equation (8), it is usually the parameter γ that is of interest because it corresponds to the DID estimand, i.e., the numerator of $\bar{\theta}_L$ in (6), and that of $\bar{\theta}_L^{(r)}$ in (7). The following theorem shows that the persuasion rates can also be obtained from the regression in (8).

Theorem 2. *Suppose that $\mathbb{E}(Y_{i1} \mid D_{i1}) > 0$ and $\Psi < 1$ so that both $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$ are well-defined. If assumption **D** holds, then*

$$\bar{\theta}_L = \frac{\gamma}{1 - \gamma_0 - \gamma_1 - \gamma_2} \quad \text{and} \quad \bar{\theta}_L^{(r)} = \frac{\gamma}{\gamma_0 + \gamma_1 + \gamma_2 + \gamma}. \quad (9)$$

The denominators in (9) reflect the fact that we need to adjust the DID estimand to estimate APRT and R-APRT. If we assume that we have a random sample $\{(Y_{i0}, Y_{i1}, G_i) : i = 1, 2, \dots, n\}$, inference based on theorem 2 can be done by combining the standard theory of OLS and the delta method.⁷

To gain further insights, we now present a seemingly different but actually equivalent approach. First, we define an auxiliary outcome variable for period 1 by $\tilde{Y}_{i1} := D_{i1} +$

⁷If the denominators are close to zero, then the delta method may provide a poor approximation in a finite sample. Using an Anderson-Rubin type statistic is an alternative robust method for inference.

$Y_{i1}(1 - D_{i1})$. That is, the auxiliary outcome \tilde{Y}_{i1} is the same as the original outcome Y_{i1} if $D_{i1} = 0$, but $\tilde{Y}_{i1} = 1$ if $D_{i1} = 1$. In words, it represents 'pseudo voters' who behave like the original voters when not exposed to the persuasive message, but will definitely vote for the party of interest when they do receive the persuasive message.

For either $A_i := \tilde{Y}_{i1} - Y_{i0}$ or $A_i := Y_{i1}D_{i1}$, consider the following moment conditions:

$$\mathbb{E}\{(Y_{i1} - Y_{i0}) - \beta_0 - \beta_1 A_i\} = 0 \quad \text{and} \quad \mathbb{E}\left[D_{i1}\{(Y_{i1} - Y_{i0}) - \beta_0 - \beta_1 A_i\}\right] = 0.$$

In other words, β_1 is a two-stage least squares (2SLS) regression coefficient of $Y_{i1} - Y_{i0}$ on A_i with using D_{i1} as an instrumental variable.

Theorem 3. *Suppose that $\mathbb{E}(Y_{i1} \mid D_{i1} = 1) > 0$ and $\Psi < 1$ so that both $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$ are well-defined. If $A_i = \tilde{Y}_{i1} - Y_{i0}$ is used, then $\beta_1 = \bar{\theta}_L$, while using $A_i = Y_{i1}D_{i1}$ leads to $\beta_1 = \bar{\theta}_L^{(r)}$.*

Therefore, both $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$ can be obtained by 2SLS regression.⁸ In order to understand the idea behind theorem 3, first focus on $\bar{\theta}_L$, i.e., APRT, where both the numerator and the denominator have a form of a DID estimand. Recall that a DID estimand can be obtained by running an OLS regression of first-differenced outcomes on D_{i1} in general. Specifically, the numerator of $\bar{\theta}_L$ in (6) can be estimated by running an OLS of $Y_{i1} - Y_{i0}$ on D_{i1} in general. Furthermore, the denominator of $\bar{\theta}_L$ also has a DID form by using \tilde{Y}_{i1} in lieu of Y_{i1} . Therefore, we know from those facts that

$$\bar{\theta}_L = \frac{\text{Cov}(Y_{i1} - Y_{i0}, D_{i1})/\mathbb{V}(D_{i1})}{\text{Cov}(\tilde{Y}_{i1} - Y_{i0}, D_{i1})/\mathbb{V}(D_{i1})} = \frac{\text{Cov}(Y_{i1} - Y_{i0}, D_{i1})}{\text{Cov}(\tilde{Y}_{i1} - Y_{i0}, D_{i1})}. \quad (10)$$

The case of R-APRT is similar. That is, the numerator of $\bar{\theta}_L^{(r)}$ is the same as that of $\bar{\theta}_L$, and its denominator is $\mathbb{E}(Y_{i1} \mid D_{i1} = 1) = \text{Cov}(Y_{i1}D_{i1}, D_{i1})/\mathbb{V}(D_{i1})$, where we use the fact that D_{i1} is binary, and hence $D_{i1}^2 = D_{i1}$.

Inference based on theorem 3 can be done within the framework of the generalized method of moments (GMM): we can even stack all the moment equations to obtain both

⁸If the denominators in equation (9) are close to zero, we will have a weak instrument problem in this formulation. We do not explore this direction though.

APRT and R-APRT by a single GMM procedure. We remark that the two-way fixed effect and GMM estimators are all algebraically equivalent.

When we have covariates, it is an easy option to simply “partial them out” by adding the covariates, potentially including their powers and interactions, to the regression. However, the algebraic equivalence between the fixed effect and GMM estimators does not hold in that case. Furthermore, it can induce contamination bias, resulting in misleading conclusions in severe cases, and therefore, it requires great caution: see e.g., [Blandhol et al. \(2022\)](#) and [Goldsmith-Pinkham et al. \(2022\)](#) among others. We can address the covariate issue more properly by e.g., estimating equation (8) locally around $X_i = x$ first, after which we deal with averaging over X_i to obtain APRT or R-APRT. It is then moving toward semiparametric estimation that uses first-step estimators conditional on X_i . We will discuss several options for semiparametric estimators in the following subsection.

4.2. Semiparametric Approaches. We will now be explicit about covariates to focus on semiparametric options. For brevity and clarity, we first provide a detailed description for the more complicated parameter APRT, and then we present a short summary for R-APRT at the end of this subsection. We can construct several semiparametric estimators of $\bar{\theta}_L$ by using the first-step estimators of $\Pi_t(d, x)$ or those of $P(x) := \mathbb{P}(D_1 = 1 \mid X = x)$.⁹

First, by using the definition of $\Psi(X)$, we rewrite $\bar{\theta}_L$ as

$$\bar{\theta}_L = \frac{\mathbb{E}[\{\Delta(1, X) - \Delta(0, X)\}D_1]}{\mathbb{E}[\{\Delta(1, X) - \Delta(0, X)\}D_1] + \mathbb{E}[\{1 - \Pi_1(1, X)\}D_1]}, \quad (11)$$

where $\Delta(d, x) := \Pi_1(d, x) - \Pi_0(d, x)$. Therefore, if a random sample $\{(Y_{i0}, Y_{i1}, D_{i1}, X_i) : i = 1, \dots, n\}$ is given, then equation (11) suggests the following DID-based estimator

$$\hat{\theta}_{L, DID} := \frac{\sum_{i=1}^n \{\hat{\Delta}(1, X_i) - \hat{\Delta}(0, X_i)\}D_{i1}}{\sum_{i=1}^n \{\hat{\Delta}(1, X_i) - \hat{\Delta}(0, X_i)\}D_{i1} + \sum_{i=1}^n \{1 - \hat{\Pi}_1(1, X_i)\}D_{i1}},$$

⁹To minimize misspecification, the first-step estimators are typically nonparametric, though they may also be parametric. We do not limit ourselves to a specific type. Strictly speaking, the parametric case should be referred to as two-step parametric estimation, but for simplicity, we use the term semiparametric estimation throughout our discussion, though it may be a misnomer.

where $\widehat{\Delta}(d, X_i) := \widehat{\Pi}_1(d, X_i) - \widehat{\Pi}_0(d, X_i)$ with $\widehat{\Pi}_t(d, x)$ being the first-step estimator of $\Pi_t(d, x)$ of the researcher's choice.

There are alternative estimators. For instance, we can equivalently express $\bar{\theta}_L$ as

$$\bar{\theta}_L = \frac{\mathbb{E}\{D_1(Y_1 - Y_0)\} - \mathbb{E}\{D_1\Delta(0, X)\}}{\mathbb{E}\{D_1(1 - Y_0)\} - \mathbb{E}\{D_1\Delta(0, X)\}}, \quad (12)$$

and therefore, a direct plug-in estimator based on (12) can be obtained by

$$\hat{\theta}_{L,PI} := \frac{\sum_{i=1}^n (Y_{i1} - Y_{i0})D_{i1} - \sum_{i=1}^n \widehat{\Delta}(0, X_i)D_{i1}}{\sum_{i=1}^n (1 - Y_{i0})D_{i1} - \sum_{i=1}^n \widehat{\Delta}(0, X_i)D_{i1}}, \quad (13)$$

which will numerically coincide with $\hat{\theta}_{L,DID}$ if there are no covariates X_i .¹⁰ Further, applying the law of iterated expectations and Bayes' rule to equation (12), $\bar{\theta}_L$ can be expressed in another form as follows:

$$\bar{\theta}_L = \frac{\mathcal{N}}{\mathcal{N} + \mathbb{E}\{D_1(1 - Y_1)\}}, \quad (14)$$

where

$$\mathcal{N} := \mathbb{E}\left\{D_1(Y_1 - Y_0) - (1 - D_1)(Y_1 - Y_0)\frac{P(X)}{1 - P(X)}\right\}. \quad (15)$$

Therefore, we can obtain yet another estimator from equation (14), i.e.,

$$\hat{\theta}_{L,POW} := \frac{\widehat{\mathcal{N}}}{\widehat{\mathcal{N}} + n^{-1} \sum_{i=1}^n (1 - Y_{i1})D_{i1}}, \quad (16)$$

where

$$\widehat{\mathcal{N}} := \frac{1}{n} \sum_{i=1}^n (Y_{i1} - Y_{i0})D_{i1} - \frac{1}{n} \sum_{i=1}^n (1 - D_{i1})(Y_{i1} - Y_{i0})\frac{\widehat{P}(X_i)}{1 - \widehat{P}(X_i)}$$

with $\widehat{P}(X)$ being the first-step estimator of the propensity score $P(X) := \mathbb{P}(D_1 = 1 \mid X)$ of the researcher's choice. The estimator $\hat{\theta}_{L,POW}$, which we will call a propensity-odds-weighted (POW) estimator, requires only one first-step estimator. Indeed, it is worth noting that $\hat{\theta}_{L,POW}$ requires estimating only $P(X)$ in the first step, while $\hat{\theta}_{L,PI}$ needs to use

¹⁰In fact, if there are no covariates, an IV estimator based on equation (10) will be identical as well.

$\Delta(0, X) = \Pi_1(0, X) - \Pi_0(0, X)$. This difference will be discussed again in the following section, where we will propose an additional estimator that is doubly (and locally) robust.

Analogously, we obtain the following estimators for R-APRT:

$$\begin{aligned}\hat{\theta}_{L,DID}^{(r)} &:= \frac{\sum_{i=1}^n \{\widehat{\Delta}(1, X_i) - \widehat{\Delta}(0, X_i)\} D_{i1}}{\sum_{i=1}^n \widehat{\Pi}_1(1, X_i) D_{i1}}, \\ \hat{\theta}_{L,PI}^{(r)} &:= \frac{\sum_{i=1}^n (Y_{i1} - Y_{i0}) D_{i1} - \sum_{i=1}^n \widehat{\Delta}(0, X_i) D_{i1}}{\sum_{i=1}^n Y_{i1} D_{i1}}, \\ \hat{\theta}_{L,POW}^{(r)} &:= \frac{\widehat{\mathcal{N}}}{n^{-1} \sum_{i=1}^n Y_{i1} D_{i1}}.\end{aligned}$$

The suggested estimators use equivalent expressions of the same moment conditions, and therefore, they will have the same efficient influence function for each parameter: i.e., they are all asymptotically equivalent under suitable regularity conditions. In the following section, we will explicitly calculate the efficient influence function for two purposes. First, it will show us the asymptotic variance of an asymptotically linear and regular semiparametric estimator of $\bar{\theta}_L$ and that of $\bar{\theta}_L^{(r)}$ without specifying all regularity conditions. Second, our derivation of the efficient influence function will be useful to find locally and doubly robust estimators.

5. THE EFFICIENT INFLUENCE FUNCTION

As it is more complicated to derive the efficient influence function for APRT, we first focus on $\bar{\theta}_L$, and we will briefly discuss the case of R-APRT at the end of this section. All the aforementioned semiparametric estimators will be asymptotically linear, regular, and normal under suitable regularity conditions: the specific conditions will depend on the choice of the first-step estimators. Since the theory of semiparametric estimation is well established (e.g., [Ackerberg et al., 2014](#)), we will not elaborate all regularity conditions to obtain asymptotic normality. Instead, we follow the approach of [Newey \(1994\)](#): i.e., we calculate the semiparametrically efficient influence function for $\bar{\theta}_L$, of which the variance will be the asymptotic variance of a regular and asymptotically linear estimator of $\bar{\theta}_L$. For this purpose, we work under the assumption of random sampling of $(Y_1, Y_0, D_1, X^\top)^\top$.

We start with making a regularity assumption.

Assumption E. *There exists a constant $\epsilon > 0$ such that for all $d, y_0, y_1 \in \{0, 1\}$ and for all $x \in \mathcal{X}$, $\epsilon \leq \mathbb{P}(Y_0 = y_0, Y_1 = y_1, D_1 = d \mid X = x) \leq 1 - \epsilon$.*

Assumption E is to ensure that the likelihood and scores are all well-behaved. Below we calculate the efficient influence function for $\bar{\theta}_L$, which we present in a couple of equivalent forms. The two forms reflect whether we estimate $P(X)$ or $\Pi_t(d, X)$'s in the first step, and they will lead us to a doubly robust estimator.

Let $\mathcal{D} := (Y_0, Y_1, D_1, X)$, and let $\theta_{L,den} := \mathbb{E}[\{1 - \Psi(X)\}D_1]$. For $EST \in \{POW, PI\}$, define

$$F_{EST,main}(\mathcal{D}) := \frac{1}{\theta_{L,den}} \left\{ H_{EST,num}(\mathcal{D}) - \theta_L H_{EST,den}(\mathcal{D}) \right\},$$

$$F_{EST,adj}(\mathcal{D}) := \frac{(1 - \theta_L)}{\theta_{L,den}} H_{EST,adj}(\mathcal{D}),$$

where

$$H_{POW,num}(\mathcal{D}) := D_1(Y_1 - Y_0) - \frac{P(X)}{1 - P(X)}(1 - D_1)(Y_1 - Y_0),$$

$$H_{POW,den}(\mathcal{D}) := D_1(1 - Y_0) - \frac{P(X)}{1 - P(X)}(1 - D_1)(Y_1 - Y_0),$$

$$H_{POW,adj}(\mathcal{D}) := -\left\{ D_1 - \frac{P(X)}{1 - P(X)}(1 - D_1) \right\} \Delta(0, X),$$

$$H_{PI,num}(\mathcal{D}) := D_1\{(Y_1 - Y_0) - \Delta(0, X)\},$$

$$H_{PI,den}(\mathcal{D}) := D_1\{(1 - Y_0) - \Delta(0, X)\},$$

$$H_{PI,adj}(\mathcal{D}) := -\frac{P(X)}{1 - P(X)}(1 - D_1)\{(Y_1 - Y_0) - \Delta(0, X)\}.$$

We are now ready to state the main theorem of this section.

Theorem 4. *Suppose that assumption E is satisfied. Then, the semiparametrically efficient influence function for $\bar{\theta}_L$ under the random sampling of $(Y_1, Y_0, D_1, X^\top)^\top$ is given by*

$$F_{DID}(\mathcal{D}) := F_{POW,main}(\mathcal{D}) + F_{POW,adj}(\mathcal{D}), \quad (17)$$

which can be equivalently written as

$$F_{DID}(\mathcal{D}) = F_{PI,main}(\mathcal{D}) + F_{PI,adj}(\mathcal{D}). \quad (18)$$

In particular, if a semiparametric estimator $\hat{\theta}_L$ of $\bar{\theta}_L$ that uses a random sample of size n is regular and asymptotically linear, then we must have

$$\sqrt{n}(\hat{\theta}_L - \bar{\theta}_L) \xrightarrow{d} N(0, \mathbb{E}\{F_{DID}^2(\mathcal{D})\}). \quad (19)$$

Theorem 4 presents the efficient influence function F_{DID} in two forms. Although equivalence of the two expressions can be easily verified by simple algebra, it is worth discussing the ideas behind them. Recall that $\hat{\theta}_{L,POW}$ and $\hat{\theta}_{L,PI}$ are alternative estimators of $\bar{\theta}_L$, where the former uses $P(X)$, while the latter does $\Delta(0, X) = \Pi_1(0, X) - \Pi_0(0, X)$. In fact, $\hat{\theta}_{L,POW}$ is based on the moment condition

$$\mathbb{E}\{H_{POW,num}(\mathcal{D}) - \bar{\theta}_L H_{POW,den}(\mathcal{D})\} = 0,$$

which corresponds to the leading term $F_{POW,main}(\mathcal{D})$ in the expression in (17). Then, $F_{POW,adj}(\mathcal{D})$ accounts for the effect of nonparametric estimation of $P(X)$. Similarly, the moment condition $\mathbb{E}\{F_{PI,main}(\mathcal{D})\} = 0$ leads to the estimator $\hat{\theta}_{L,PI}$, and the adjustment term $F_{PI,adj}(\mathcal{D})$ reflects that $\Delta(0, X) = \Pi_1(0, X) - \Pi_0(0, X)$ is estimated in the first step.

The asymptotic variance in (19) is a consequence of Theorem 2.1 in Newey (1994): because the set of scores is sufficiently rich, all regular and asymptotically linear estimators of $\bar{\theta}_L$ must have the same efficient influence function F_{DID} . An implication is that $\hat{\theta}_{L,DID}$, $\hat{\theta}_{L,PI}$, and $\hat{\theta}_{L,POW}$ will be asymptotically equivalent, as long as they are regular and asymptotically linear.

The efficient influence function formula is also useful to find a doubly robust estimator. Instead of using only the leading term in either (17) or (18), we may use the entire $F_{DID}(\mathcal{D})$ to find an estimator of $\bar{\theta}_L$. Indeed, since we trivially have

$$H_{POW,R}(\mathcal{D}) + H_{POW,adj}(\mathcal{D}) = H_{PI,R}(\mathcal{D}) + H_{PI,adj}(\mathcal{D}) \text{ for } R \in \{num, den\}, \quad (20)$$

using the moment condition $\mathbb{E}\{F_{DID}(\mathcal{D})\} = 0$, regardless of whether we use (17) or (18), leads to the same estimator, i.e.,

$$\hat{\theta}_{L,DR} := \frac{\sum_{i=1}^n (Y_{i1} - Y_{i0})D_{i1} - \sum_{i=1}^n \hat{\Delta}(0, X_i)D_{i1} + \sum_{i=1}^n \hat{H}_{PI,adj}(\mathcal{D}_i)}{\sum_{i=1}^n (1 - Y_{i0})D_{i1} - \sum_{i=1}^n \hat{\Delta}(0, X_i)D_{i1} + \sum_{i=1}^n \hat{H}_{PI,adj}(\mathcal{D}_i)}, \quad (21)$$

where

$$\hat{H}_{PI,adj}(\mathcal{D}_i) := -\frac{\hat{P}(X_i)}{1 - \hat{P}(X_i)}(1 - D_{i1})\{(Y_{i1} - Y_{i0}) - \hat{\Delta}(0, X_i)\}.$$

By construction, the estimator $\hat{\theta}_{L,DR}$ is both locally and doubly robust. Specifically, since it is using the entire efficient influence function, the effect of first-step nonparametric estimation is already reflected in the form of the estimator: the asymptotic variance of $\hat{\theta}_{L,DR}$ will be the same as its oracle form that uses the true $\Delta(0, X)$ and $P(X)$. This local robustness suggests that we can develop a double/debiased machine learning (DML) estimator in a similar way to $\hat{\theta}_{L,DR}$ but by using machine learning estimators in the first step and cross-fitting in the second step. Further, in view of equation (20), we can see that $\hat{\theta}_{L,DR}$ is doubly robust in that it will be consistent as long as either $\hat{\Delta}(0, \cdot)$ or $\hat{P}(\cdot)$, but not necessarily both, is correctly specified: i.e., $\mathbb{E}\{H_{PI,adj}(\mathcal{D})\} = 0$ holds even if $P(X)$ is replaced with an arbitrary function of X , while $\mathbb{E}\{H_{POW,adj}(\mathcal{D})\} = 0$ stays true whatever function of X is used in lieu of $\Delta(0, X)$.

We end this section by discussing the case of R-APRT. The only difference between APRT and R-APRT is in their denominators. Indeed, R-APRT has a simpler denominator, i.e., $\theta_{L,den}^{(r)} := \mathbb{E}(Y_1 D_1)$, and it can be shown that under the same condition as theorem 4, the efficient influence function for $\theta_L^{(r)}$ is given by

$$F_{DID}^{(r)}(\mathcal{D}) := \frac{1}{\theta_{L,den}^{(r)}} \left\{ H_{EST,num}(\mathcal{D}) - \theta_L^{(r)} Y_1 D_1 + H_{EST,adj}(\mathcal{D}) \right\},$$

where $EST \in \{POW, PI\}$. Therefore, the same reasoning as before shows that a (locally and) doubly robust estimator of $\theta_L^{(r)}$ can be obtained by

$$\hat{\theta}_{L,DR}^{(r)} := \frac{\sum_{i=1}^n (Y_{i1} - Y_{i0})D_{i1} - \sum_{i=1}^n \hat{\Delta}(0, X_i)D_{i1} + \sum_{i=1}^n \hat{H}_{PI,adj}(\mathcal{D}_i)}{\sum_{i=1}^n Y_{i1} D_{i1}}.$$

6. DISCUSSION

In this section, we discuss how to conduct back-of-the-envelope inference on APRT and R-APRT when we do not have access to the full data, but we have information for ATT. For the sake of simple discussion, we again focus on the two time-period case and assume that ATT is weakly positive (as implied by assumption B). APRT is linked to ATT via $\text{APRT}(q) = \text{ATT}/(\text{ATT} + q)$, where $q := \mathbb{P}(Y_1 = 0 \mid D_1 = 1)$. Because $\text{ATT} + q = \mathbb{P}\{Y_1(0) = 0 \mid D_1 = 1\}$, APRT is strictly larger than ATT unless $\text{ATT} = 0$ or $\mathbb{P}\{Y_1(0) = 0 \mid D_1 = 1\} = 1$, which means that every individual in the treated group is a target audience to persuade. In other words, APRT will be strictly larger than ATT in many interesting cases.

Suppose that there is a known interval $[\underline{q}, \bar{q}] \subseteq [0, 1]$ such that it contains q with probability tending to $1 - \alpha_0$ for some constant $0 \leq \alpha_0 < 1$. As $q \mapsto \text{APRT}(q)$ is nonincreasing in q (under the assumption that $\text{ATT} \geq 0$), the resulting bounds on APRT conditional on $q \in [\underline{q}, \bar{q}]$ are

$$[L(\text{APRT}), U(\text{APRT})] := \left[\frac{\text{ATT}}{\text{ATT} + \bar{q}}, \frac{\text{ATT}}{\text{ATT} + \underline{q}} \right], \quad (22)$$

which is our basis for the back-of-the-envelope inference on APRT from ATT.

To be more specific, suppose that an estimate $\widehat{\text{ATT}}$ of ATT, and its standard error, say $se(\widehat{\text{ATT}})$, are available so that the asymptotic $(1 - \alpha)$ confidence interval for ATT is obtained by $[\widehat{\text{ATT}} - z_{1-\alpha/2} \cdot se(\widehat{\text{ATT}}), \widehat{\text{ATT}} + z_{1-\alpha/2} \cdot se(\widehat{\text{ATT}})]$, where z_τ is the τ quantile of the standard normal distribution. Then, by the delta method, the pointwise standard error of $\widehat{\text{APRT}}(q) := \widehat{\text{ATT}}/(\widehat{\text{ATT}} + q)$ is $se(\widehat{\text{ATT}})q/(\widehat{\text{ATT}} + q)^2$. Then, in view of the interval identification in (22), a $(1 - \alpha)$ Bonferroni confidence interval for APRT can be obtained by $[\underline{\text{APRT}}, \overline{\text{APRT}}]$, where

$$\begin{aligned} \underline{\text{APRT}} &:= \widehat{\text{APRT}}(\bar{q}) - z_{1-(\alpha-\alpha_0)/2} \cdot \frac{se(\widehat{\text{ATT}})\bar{q}}{(\widehat{\text{ATT}} + \bar{q})^2}, \\ \overline{\text{APRT}} &:= \widehat{\text{APRT}}(\underline{q}) + z_{1-(\alpha-\alpha_0)/2} \cdot \frac{se(\widehat{\text{ATT}})\underline{q}}{(\widehat{\text{ATT}} + \underline{q})^2}. \end{aligned}$$

because

$$\begin{aligned} \mathbb{P}\{\text{APRT} \in [\underline{\text{APRT}}, \overline{\text{APRT}}]\} &\geq \mathbb{P}\{q \in [\underline{q}, \bar{q}], \text{APRT} \in [\underline{\text{APRT}}, \overline{\text{APRT}}]\} \\ &\geq 1 - \mathbb{P}(q \notin [\underline{q}, \bar{q}]) - \mathbb{P}(\text{APRT} \notin [\underline{\text{APRT}}, \overline{\text{APRT}}]) = 1 - \alpha_0 - \frac{\alpha - \alpha_0}{2} - \frac{\alpha - \alpha_0}{2} = 1 - \alpha. \end{aligned}$$

For example, we may set $\alpha_0 = \alpha/2$, resulting in $(\alpha - \alpha_0)/2 = \alpha/4$. In short, when we have access to the estimate of ATT with its standard error, we can conduct back-of-the-envelope inference on APRT, provided that we have probabilistic bounds on $\mathbb{P}(Y_1 = 0 \mid D_1 = 1)$.

We now move to R-APRT. Using the fact that $\text{R-APRT} = \text{ATT}/(1 - q)$, we propose the following back-of-the-envelope confidence interval for R-APRT:

$$\left[\widehat{\text{R-APRT}}(\underline{q}) - z_{1-(\alpha-\alpha_0)/2} \cdot \frac{se(\widehat{\text{ATT}})}{(1-\underline{q})}, \widehat{\text{R-APRT}}(\bar{q}) - z_{1-(\alpha-\alpha_0)/2} \cdot \frac{se(\widehat{\text{ATT}})}{(1-\bar{q})} \right].$$

7. EXAMPLE I: NEWS MEDIA PERSUASION

Ladd and Lenz (2009) exploited abrupt shifts in British newspaper endorsements from the Conservative party to the Labour party before the 1997 general election. Using data from the British Election Panel Study, they compared readers of newspapers that switched endorsements (treated group) with those who did not read these newspapers (control group). The binary outcome variable is whether a respondent voted Labour in the 1997 election. The binary treatment variable is whether a respondent read the switching newspapers. This is an example of the two-period case with no treatment in the first period. In addition to the lagged outcome (prior Labour vote in 1992), there are a large number of pretreatment variables (X) all measured in 1992. The dataset is on the public domain as part of replication materials for **Hainmueller (2012)**, which are available in the Political Analysis Dataverse at <http://dvn.iq.harvard.edu/dvn/dv/pan>. We use the same set of covariates as in **Hainmueller (2012)** but exclude the three prior voting variables (prior Labour vote, prior Conservative vote, prior Liberal vote), resulting in 36 predetermined covariates. They include prior party identification, prior party support, prior ideology,

parents vote, political knowledge, television viewership, daily newspaper readership, authoritarianism, trade union membership, mortgage payment status, education, income, age, race, socioeconomic status, gender, region, and occupation.

Table 1 reports a variety of estimates of the average persuasion rates on the treated. Panel A reports the forward version, i.e., APRT and panel B the reverse version, i.e., R-APRT. In the columns of FE and GMM, the two-way fixed effect (FE) and GMM estimators are given without controlling for the covariates. They are algebraically equivalent but the t-statistics are slightly different because they rely on different asymptotic approximations. The covariates are separately partialled out in the columns of FE-X and GMM-X. Each of the two-step estimates, which are given in the columns of DID, PI, POW, and DR, involves the first step estimation of (some but not all of) five conditional probabilities: $\mathbb{P}(Y_t = 1 \mid D_1 = d, X = x)$, where $t = 0, 1$ and $d = 0, 1$, as well as $\mathbb{P}(D_1 = 1 \mid X = x)$. They are all estimated via logistic regression that is linear in X .

TABLE 1. Average Persuasion Rates on the Treated: Media Persuasion

Panel A: APRT								
method	Regression-Based				Two-Step			
	FE	GMM	FE-X	GMM-X	DID	PI	POW	DR
equation	(9)	(10)	(9)	(10)	(11)	(12)	(14)	(17)
covariates	none		partial out		two-step estimation			
estimate	0.172	0.172	0.170	0.169	0.177	0.172	0.184	0.176
t-statistic	2.798	2.800	2.384	2.777	2.687	2.609	2.792	2.667
Panel B: R-APRT								
estimate	0.148	0.148	0.150	0.149	0.149	0.149	0.162	0.153
t-statistic	2.737	2.739	2.522	2.761	2.651	2.645	2.872	2.716

To compare our estimates with other measures reported in the literature, we note that our ATT estimate using the doubly robust method is just 0.089, which is very similar to the unconditional DID estimate of 0.086 (Ladd and Lenz, 2009, see the first column of their Table 2). The previous ATT estimates using matching (Ladd and Lenz, 2009) or entropy balancing (Hainmueller, 2012) are slightly larger, ranging from 0.096 to 0.140. Our

estimates of APRT and R-APRT are larger than any of the ATT estimates, indicating that media persuasion was even more substantial than the ATT estimates suggest. Using the decomposition method described in our discussion below equation (2) and in the one that follows assumption B, we can divide the population of the treated by three groups: the share of the treatment-persuadable is 0.089 as it is just ATT; that of the never-persuadable is 0.417, and that of the already-persuaded is 0.494. The shares of the latter two groups are much larger, explaining why ATT is substantially different from APRT and R-APRT.

We end this section by illustrating back-of-the-envelope inference. Strictly speaking, we do not have access to the replication files of the original analysis in Ladd and Lenz (2009) because we rely on the replication materials from Hainmueller (2012). Table 2 in Ladd and Lenz (2009) provides a variety of treatment effect estimates. For example, the third column shows the DID estimate using only exactly matched treated/control groups, where exact matching is implemented on selected variables. The resulting estimate is 0.109 with its standard error of 0.041. We can interpret this as the ATT estimate. In order to conduct inference on APRT using their ATT estimate, we need bounds on $q = \mathbb{P}(Y_1 = 0 \mid D_1 = 1)$. Its estimated value, 0.583, and the size, 211, of the treated group are given on page 399 in their paper. To be conservative, we use the 97.5% confidence interval on q , that is $[q, \bar{q}] = [0.507, 0.659]$. Then, the 95% confidence interval on APRT is $[0.039, 0.300]$, while the point estimate is 0.158 using $q = 0.583$. If we carry out the same exercise for R-APRT, we obtain the back-of-the-envelope confidence interval of $[0.035, 0.589]$ with the point estimate of R-APRT being 0.261.

8. STAGGERED TREATMENT

As we extend our setup to allow for multiple time periods, we focus on the leading case of staggered treatments, where the observational units receive treatment at different times and the treatment stays on once it is adopted. Furthermore, to avoid repetition, we limit our attention to the more complicated parameter, i.e., APRT, noting that the modifications required for R-APRT are easily adaptable.

8.1. The Setup and the Parameters of Interest. Let $t \in \{0, 1, 2, \dots, T\}$ denote the time period, and let D_t denote the treatment status at time t . In our setup, no one is treated at time 0, and if the agent receives a treatment at time $s \in \{1, 2, \dots, T\}$, then we observe $D_t = 0$ for all $0 \leq t < s$ and $D_t = 1$ for all $s \leq t \leq T$. If the agents are never treated so that $D_t = 0$ for all $t \in \{0, 1, 2, \dots, T\}$, then they belong to the control group. Exogenous covariates are denoted by X as before.

Let $\mathcal{S} := \{t \in \{1, \dots, T\} : D_t = 1\}$, and define the new random variable S by $S := \min \mathcal{S}$ if \mathcal{S} is non-empty, and $S := \infty$ otherwise. Therefore, S has the support $\{1, 2, \dots, T, \infty\}$, indicating the time period at which the treatment is given: i.e., $S = \infty$ means that the individuals are never treated.

Now, in order to capture treatment effect heterogeneity, we write $Y_t(s)$ for the potential outcome at time t for the case where the treatment is given at time s . We now make the following assumption to formally describe our setting.

Assumption F. *At time 0, no one receives or anticipates receiving a treatment so that $Y_0 := Y_0(\infty)$ is observed. At time $t = 1, \dots, T$, $Y_t := Y_t(S)$ is observed, and there is no anticipation for a treatment in that we have $\mathbb{P}\{Y_t(s) = 1 \mid S = s, X\} = \mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X\}$ whenever $t < s$ with probability one. Further, for all $s, t \in \{1, \dots, T\}$ and for (almost) all $x \in \mathcal{X}$, there is a constant $\epsilon > 0$ such that $\epsilon \leq \mathbb{P}\{Y_t(\infty) = 1, S = s \mid X = x\}$ and $\mathbb{P}(S = s \mid X = x) \leq 1 - \epsilon$.*

The most straightforward parameters we can start with are the persuasion rate at time t on those who are treated at time s conditional on $X = x$, and its aggregation across X :

$$\theta_c(s, t \mid x) := \mathbb{P}\{Y_t(s) = 1 \mid Y_t(\infty) = 0, S = s, X = x\},$$

$$\theta(s, t) := \mathbb{P}\{Y_t(s) = 1 \mid Y_t(\infty) = 0, S = s\}.$$

Interpretation is not any different from that of the persuasion rate on the treated in the case of two periods. For example, $\theta(s, t)$ measures the proportion of those who take the action of interest at time t (e.g. voting for a certain party) among those who received a persuasive message at time s and who would not have taken the action without the

persuasive message at all. Again, the point is that we focus only on those who switch their action because of the persuasive message.

As a measure of the cumulative effect of persuasion, we consider the j -period-forward persuasion rate on the treated, i.e. the persuasion rate in j periods after the time of the treatment on those who are ever treated (by $T - j$). Following the literature on the panel event-study design, we call it *the event-study persuasion rate* (ESPR), and it can be formally expressed by

$$\theta_{ESPR}(j) := \mathbb{P}\{Y_{S+j}(S) = 1 \mid Y_{S+j}(\infty) = 0, S \leq T - j\} \text{ for } j = 0, 1, \dots, T - 1.$$

This is probably the most comprehensive summary parameter, and its identification, estimation, and inference will be discussed in detail in the following subsections.¹¹

8.2. Identification. All the parameters defined above depend on the joint probabilities of the potential outcomes. As in the case of two periods, a monotonicity assumption is useful to handle the situation.

Assumption G. For all $s, t \in \{1, \dots, T\}$, $\mathbb{P}\{Y_t(s) \geq Y_t(\infty) \mid S = s, X\} = 1$ almost surely.

Assumption G is an extension of assumption B to the case of staggered treatment: i.e., the persuasive message is directional, and there is no backlash regardless of the timing of the exposure, and therefore the monotonicity condition is satisfied.

As before, all the parameters will be identified if

$$\tau_{stagger}(s, t \mid x) := \mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X = x\} \quad \text{with } t \geq s$$

is identified for all $x \in \mathcal{X}$. Identification of $\tau_{stagger}(s, t \mid x)$ can be achieved via the DID approach. For this purpose, we modify assumption C as follows.

¹¹Alternatively, $\theta(s, t)$ could be aggregated in different ways; for instance, a new summary parameter could be constructed as a function of t .

Assumption H. For all $s, t \in \{0, 1, 2, \dots, T\}$ and for (almost) all $x \in \mathcal{X}$, $\mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X = x\}$ is separable into the sum of a time component and a treatment component: i.e., there exist functions G^* and H^* such that $\mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X = x\} = G^*(t, x) + H^*(s, x)$.

Assumption **H** is an extended version of the parallel trend assumption. Specifically, assumption **H** ensures that for any $s, t \in \{1, 2, \dots, T\}$ with $t \geq s$, we have

$$\begin{aligned} & \mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X = x\} - \mathbb{P}\{Y_{s-1}(\infty) = 1 \mid S = s, X = x\} \\ &= \mathbb{P}\{Y_t(\infty) = 1 \mid S = \infty, X = x\} - \mathbb{P}\{Y_{s-1}(\infty) = 1 \mid S = \infty, X = x\}. \end{aligned} \quad (23)$$

Let

$$\begin{aligned} \Psi_{stagger}(s, t \mid x) &:= \mathbb{P}(Y_{s-1} = 1 \mid S = s, X = x) \\ &+ \mathbb{P}(Y_t = 1 \mid S = \infty, X = x) - \mathbb{P}(Y_{s-1} = 1 \mid S = \infty, X = x), \end{aligned}$$

which is directly identified from panel data. We have the following identification results.

Theorem 5. Suppose that assumptions **F** to **H** are satisfied. Then, $\theta_c(s, t \mid x)$, $\theta(s, t)$, and $\theta_{ESPR}(j)$ are point-identified by

$$\theta_c(s, t \mid x) = \frac{\mathbb{P}(Y_t = 1 \mid S = s, X = x) - \Psi_{stagger}(s, t \mid x)}{1 - \Psi_{stagger}(s, t \mid x)}, \quad (24)$$

$$\theta(s, t) = \frac{\mathbb{P}(Y_t = 1 \mid S = s) - \mathbb{E}\{\Psi_{stagger}(s, t \mid X) \mid S = s\}}{1 - \mathbb{E}\{\Psi_{stagger}(s, t \mid X) \mid S = s\}}, \quad (25)$$

$$\theta_{ESPR}(j) = \frac{\sum_{s=1}^{T-j} \mathbb{P}(S = s) [\mathbb{P}(Y_{s+j} = 1 \mid S = s) - \mathbb{E}\{\Psi_{stagger}(s, s+j \mid X) \mid S = s\}]}{\sum_{s=1}^{T-j} \mathbb{P}(S = s) [1 - \mathbb{E}\{\Psi_{stagger}(s, s+j \mid X) \mid S = s\}]}. \quad (26)$$

Theorem **5** follows from the intermediate identification result that under assumption **H**, $\tau_{stagger}(s, t \mid x) = \Psi_{stagger}(s, t \mid x)$. The fact that the numerator and the denominator of equation (24) are separately aggregated to obtain equation (25) is a consequence of the Bayes rule. Further, it is worth noting that ESPR is obtained by aggregating over s on the numerator and on the denominator separately, which is reminiscent of the aggregation procedure for X . Since the persuasion rate is a conditional probability, it takes the form of

a ratio, but when it comes to aggregation (over X or s), it is always collapsed to separate aggregations of the numerator and denominator.

8.3. Estimation. All the estimators we discussed in the case of two periods can be extended to the current setup of staggered treatment. Below we discuss both regression-based and semiparametric approaches, focusing on $\theta_{ESPR}(j)$ as a comprehensive summary parameter.

8.3.1. Regression-Based Approaches. We assume that we have no covariates for now. Theorem 5 implies that

$$\theta(s, s + j) = \frac{\theta_{num}(s, s + j)}{\theta_{den}(s, s + j)} \quad \text{and} \quad \theta_{ESPR}(j) = \frac{\sum_{s=1}^{T-j} \theta_{num}(s, s + j) \mathbb{P}(S = s)}{\sum_{s=1}^{T-j} \theta_{den}(s, s + j) \mathbb{P}(S = s)},$$

where

$$\begin{aligned} \theta_{num}(s, s + j) &:= \mathbb{E}(Y_{s+j} | S = s) - \mathbb{E}(Y_{s-1} | S = s) - \mathbb{E}(Y_{s+j} | S = \infty) + \mathbb{E}(Y_{s-1} | S = \infty), \\ \theta_{den}(s, s + j) &:= 1 - \mathbb{E}(Y_{s-1} | S = s) - \mathbb{E}(Y_{s+j} | S = \infty) + \mathbb{E}(Y_{s-1} | S = \infty). \end{aligned}$$

Below we use linear regression to estimate $\theta(s, s + j)$, which has a similar structure to the persuasion rate in the two-period case. For instance, the numerator $\theta_{num}(s, s + j)$ is a DID estimand, where $S = s$ represents the treatment group, and $S = \infty$ the control group. Therefore, we can extend the ideas in section 4.1 to the current setup in a pairwise manner.

To be more specific, suppose that we have a random sample $\{(Y_{i0}, Y_{i1}, \dots, Y_{iT}, S_i) : i = 1, \dots, n\}$. For $\theta(s, s + j)$, we use only the subset of the data $\{(Y_{is-1}, Y_{s+j}, S_i) : S_i = s \text{ or } S_i = \infty\}$, and we estimate the following regression by OLS:

$$Y_{it} = \gamma_0 + \mathbb{1}(S_i = s)\gamma_1 + \mathbb{1}(t = s + j)\gamma_2 + \mathbb{1}(S_i = s)\mathbb{1}(t = s + j)\gamma + \epsilon_{it}, \quad (27)$$

where we assume that $\mathbb{E}(\epsilon_{it} \mid S_i = d, t = r) = 0$ for $d \in \{s, \infty\}$ and $r \in \{s - 1, s + j\}$. The extra conditioning of $S_i \in \{s, \infty\}$ is to ensure that $S_i \neq s$ corresponds to $S_i = \infty$ and vice versa.¹² It is in the same spirit to restrict ourselves to $t \in \{s - 1, s + j\}$.

Equation (27) can be interpreted like the usual two-way fixed effect regression: γ_1 captures the group effect, γ_2 does the time effect, and γ is the DID estimand that corresponds to the numerator $\theta_{num}(s, s + j)$. The denominator $\theta_{den}(s, s + j)$ can be obtained from the rest of the coefficients as we did in theorem 2.

Theorem 6. *Consider the regression of (27). If we restrict ourselves to the data with $S_i \in \{s, \infty\}$ and $t \in \{s - 1, s + j\}$, then we have*

$$\theta_{num}(s, s + j) = \gamma \quad \text{and} \quad \theta_{den}(s, s + j) = 1 - \gamma_0 - \gamma_1 - \gamma_2.$$

Theorem 6 suggests the following sample analog estimators:

$$\hat{\theta}(s, s + j) := \frac{\hat{\theta}_{num}(s, s + j)}{\hat{\theta}_{den}(s, s + j)} \quad \text{and} \quad \hat{\theta}_{ESPR}(j) := \frac{\sum_{s=1}^{T-j} \hat{\theta}_{num}(s, s + j) \hat{\mathbb{P}}(S = s)}{\sum_{s=1}^{T-j} \hat{\theta}_{den}(s, s + j) \hat{\mathbb{P}}(S = s)}, \quad (28)$$

where $\hat{\theta}_{num}(s, s + j) := \hat{\gamma}$, $\hat{\theta}_{den}(s, s + j) := 1 - \hat{\gamma}_0 - \hat{\gamma}_1 - \hat{\gamma}_2$, $(\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_2, \hat{\gamma})$ are the OLS estimators from (27), and $\hat{\mathbb{P}}(S = s)$ is the sample analog of $\mathbb{P}(S = s)$.¹³

Instead of running the regression of (27) for each pair of $(s - 1, s + j)$, one may run the following regression model that resembles the panel event-study design regression more closely: for each $s \in S$, using only the subset of the data $\{(Y_{it}, S_i) : S_i = s, t = 0, \dots, T\}$, we run the following regression:

$$y_{it} = \mu_{-1}^{(s)} + \sum_{j \neq -1} \alpha_j^{(s)} \mathbb{1}(t = s + j) + \epsilon_{it}. \quad (29)$$

¹² It is useful to recall that for any events E_1, E_2, E_3 , $\mathbb{P}(E_1 \mid E_2) = \mathbb{P}(E_1 \mid E_2, E_2 \cup E_3)$, provided that $\mathbb{P}(E_2) > 0$.

¹³ The OLS estimators $(\hat{\gamma}_0, \hat{\gamma}_1, \hat{\gamma}_2, \hat{\gamma})$ should be indexed by $(s, s + j)$, but its dependence is suppressed to minimize the notional burden.

Then, the estimators of $\theta(s, s + j)$ and $\hat{\theta}_{ESPR}(j)$ are now given by

$$\hat{\theta}(s, s + j) = \frac{\hat{\alpha}_j^{(s)} - \hat{\alpha}_j^{(\infty)}}{1 - \hat{\mu}_{-1}^{(s)} - \hat{\alpha}_j^{(\infty)}} \quad \text{and} \quad \hat{\theta}_{ESPR}(j) = \frac{\sum_{s=1}^{T-j} \{\hat{\alpha}_j^{(s)} - \hat{\alpha}_j^{(\infty)}\} \widehat{\mathbb{P}}(S = s)}{\sum_{s=1}^{T-j} \{1 - \hat{\mu}_{-1}^{(s)} - \hat{\alpha}_j^{(\infty)}\} \widehat{\mathbb{P}}(S = s)}, \quad (30)$$

where a hat on the right-hand side of the equations above denotes the OLS estimator. The estimators defined in equations (28) and (30) look seemingly different; however, they are identical because we run fully saturated regressions in two different ways. Finally, we comment that statistical inference is straightforward. Since the OLS-based estimators are all asymptotically linear, $\hat{\theta}(s, s + j)$ and $\hat{\theta}_{ESPR}(j)$ are asymptotically linear and normal by the delta method.

8.3.2. Semiparametric Approaches. We now consider the covariates and describe how to implement semiparametric estimation. First, let $\bar{S}_s = \mathbb{1}(S = s) + \mathbb{1}(S = \infty)$ for each $s \in \{1, 2, \dots, T - j\}$. Using the fact in footnote 12, we can then equivalently express $\theta(s, s + j)$ as follows:

$$\theta(s, s + j) = \frac{\mathcal{N}_{stagger}(s, s + j \mid \bar{S}_s = 1)}{\mathcal{N}_{stagger}(s, s + j \mid \bar{S}_s = 1) + \mathbb{P}(Y_{s+j} = 0, S = s \mid \bar{S}_s = 1)}, \quad (31)$$

where

$$\begin{aligned} \mathcal{N}_{stagger}(s, s + j \mid \bar{S}_s = 1) &:= \mathbb{P}(Y_{s+j} = 1, S = s \mid \bar{S}_s = 1) - \mathbb{P}(Y_{s-1} = 1, S = s \mid \bar{S}_s = 1) \\ &\quad - \mathbb{E} \left\{ (Y_{s+j} - Y_{s-1}) \mathbb{1}(S = \infty) \frac{\mathbb{P}(S = s \mid X, \bar{S}_s = 1)}{1 - \mathbb{P}(S = s \mid X, \bar{S}_s = 1)} \mid \bar{S}_s = 1 \right\}. \end{aligned}$$

These expressions conveniently show that our discussions on estimation and inference in the case of two periods can be applied to the current case of multiple periods.

For example, we can obtain a doubly robust estimator of $\theta(s, s + j)$ based on a random sample $\{(Y_{i0}, Y_{i1}, \dots, Y_{iT}, S_i, X_i^\top)^\top : i = 1, 2, \dots, n\}$ by using the same formula for the two-period DR estimator $\hat{\theta}_{L,DR}$ defined in equation (21) but using a subsample of $(Y_{is-1}, Y_{is+j}, \mathbb{1}(S_i = s), X_i^\top)^\top$ that satisfies $\mathbb{1}(S_i = s \text{ or } S_i = \infty) = 1$. Since $\mathbb{1}(S_i = s) \mathbb{1}(\bar{S}_{s,i} =$

1) = $\mathbb{1}(S_i = s)$, we obtain the following formula:

$$\hat{\theta}_{DR}(s, s + j) := \frac{\sum_{i=1}^n \{(Y_{is+j} - Y_{is-1} - \hat{\Delta}_{s,s+j}(\infty, X_i))\} \mathbb{1}(S_i = s) - \hat{C}(s, s + j)}{\sum_{i=1}^n \{(1 - Y_{is-1} - \hat{\Delta}_{s,s+j}(\infty, X_i))\} \mathbb{1}(S_i = s) - \hat{C}(s, s + j)},$$

where

$$\hat{\Delta}_{s,s+j}(\infty, X_i) := \hat{\mathbb{P}}(Y_{is+j} = 1 \mid S_i = \infty, X_i) - \hat{\mathbb{P}}(Y_{is-1} = 1 \mid S_i = \infty, X_i),$$

$$\hat{C}(s, s + j) := \sum_{i=1}^n \frac{\hat{\mathbb{P}}(S_i = s \mid X_i, \bar{S}_{s,i} = 1) \mathbb{1}(S_i = \infty)}{1 - \hat{\mathbb{P}}(S_i = s \mid X_i, \bar{S}_{s,i} = 1)} \{(Y_{is+j} - Y_{is-1}) - \hat{\Delta}_{s,s+j}(\infty, X_i)\}.$$

Aggregating over s in $\hat{\theta}_{DR}(s, s + j)$ can be similarly done. For example, a DR estimator of $\theta_{ESPR}(j)$ can be obtained by

$$\hat{\theta}_{ESPR,DR}(j) := \frac{\sum_{s=1}^{T-j} \hat{\mathbb{P}}(S = s) \left[\sum_{i=1}^n \{(Y_{is+j} - Y_{is-1} - \hat{\Delta}_{s,s+j}(\infty, X_i))\} \mathbb{1}(S_i = s) - \hat{C}(s, s + j) \right]}{\sum_{s=1}^{T-j} \hat{\mathbb{P}}(S = s) \left[\sum_{i=1}^n \{(1 - Y_{is-1} - \hat{\Delta}_{s,s+j}(\infty, X_i))\} \mathbb{1}(S_i = s) - \hat{C}(s, s + j) \right]}.$$

We discuss inference for $\hat{\theta}_{ESPR,DR}(j)$ in Online Appendix S-6.

9. EMPIRICAL EXAMPLE II: CURRICULUM PERSUASION

[Cantoni et al. \(2017\)](#) investigated the impact of a textbook reform in China by exploiting a staggered introduction of the new curriculum across provinces. Specifically, they conducted a survey of the political attitudes and beliefs with Peking University undergraduate students of four cohorts of students, who entered high school between 2006 and 2009. They adopted a difference-in-differences framework, and their baseline regression model ([Cantoni et al., 2017](#), see equation (1) on p.361) is of the following form: using their notation,

$$y_{icp} = \sum_c \gamma_c + \sum_p \delta_p + \beta \text{New Curriculum}_{cp} + \varepsilon_{icp}, \quad (32)$$

where i denotes the individual, c the high school entry cohort, and p the province of high school attendance; y_{icp} is an outcome variable, $\text{New Curriculum}_{cp}$ is the indicator variable

that has value 1 if cohort c in province p studied under the new curriculum and 0 otherwise, and γ_c and δ_p are cohort and province fixed effects; and β is the parameter of interest, aiming to uncover the effect of the new curriculum.

Based on the regression estimates of (32), they also reported the persuasion rates (Cantoni et al., 2017, see column (7) of Table 3). To describe their estimates of the persuasion rates, let $\hat{\beta}$, $\hat{\gamma}_c$ and $\hat{\delta}_p$, respectively, denote the estimates of β , γ_c , and δ_p in the two-way fixed effect regression of (32), where if y_{icp} is not binary, the outcome variable is transformed to a binary variable that equals one if the outcome is greater than or equal to the median answer. Then, the persuasion rate reported in Cantoni et al. (2017) can be written as

$$\hat{\theta}_{CCYYZ} := \frac{\hat{\beta}}{1 - (\sum_c \hat{\gamma}_c + \sum_p \hat{\delta}_p)},$$

where the numerator $\hat{\beta}$ aims to measure the impacts of the curriculum reform, while the denominator is a proxy of “the share of students without the desired attitude among individuals who studied under the old curriculum” using the words of Cantoni et al. (2017, p. 382).

If there were only two cohorts and two provinces, $\hat{\theta}_{CCYYZ}$ would be the same as our two-period regression-based estimator of APRT developed in section 4.1. In light of this, we do not consider R-APRT but focus on APRT. However, their setting involves a staggered introduction of the new curriculum with more than two cohorts and in more than two provinces. In addition, in view of recent advances in the literature on staggered treatment (e.g. De Chaisemartin and d’Haultfoeuille, 2020; Callaway and Sant’Anna, 2021), which postdated Cantoni et al. (2017), note that $\hat{\beta}$ (the numerator of $\hat{\theta}_{CCYYZ}$) might be a biased estimate of ATT. Hence, we revisit their analysis. Their sample is not from panel data but from a survey with multiple cohorts. However, a simple modification of the regression-based approach in section 8.3.1 provides a method for estimating the persuasion rates in this example. In their sample, there are four high school entry cohorts: 2006, 2007, 2008, and 2009. As we need one comparison year before treatment, the earliest treatment group

($S = 1$ using our notation) is a subset of cohort 2007 who attended high school in the provinces that introduced the new curriculum in 2007. The control group ($S = \infty$ using our notation) consists of those who attended high school in the provinces that introduced the new curriculum in 2010, which is the last year when the textbook reform was complete across China. Table 2 summarizes the staggered adoption of treatment in this example.

TABLE 2. Staggered Introduction of the New Curriculum

Group (S) by year of a new curriculum	Event Horizon (j)						
	-4	-3	-2	-1	0	1	2
2007				2006	2007	2008	2009
2008			2006	2007	2008	2009	
2009		2006	2007	2008	2009		
2010 ($S = \infty$)	2006	2007	2008	2009			

^a The provinces that introduced the new curriculum in 2007 are Beijing, Heilongjiang, Hunan, Jilin, and Shaanxi.

^b The provinces that introduced the new curriculum in 2008 are Henan, Jiangxi, Shanxi, and Xinjiang.

^c The provinces that introduced the new curriculum in 2009 are Hebei, Hubei, Inner Mongolia, and Yunnan.

^d The provinces that introduced the new curriculum in 2010 are Chongqing, Gansu, Guangxi, Guizhou, Qinghai, Sichuan, and Tibet.

^e An entry in the table refers to the relevant high school entry year (i.e., $s + j$) for each $S = s$ group.

^f The provinces that adopted the new curriculum before 2007 are excluded in our analysis.

We now describe how to estimate $\theta(s, s + j)$. As in (29), for each s , we run the following regression:

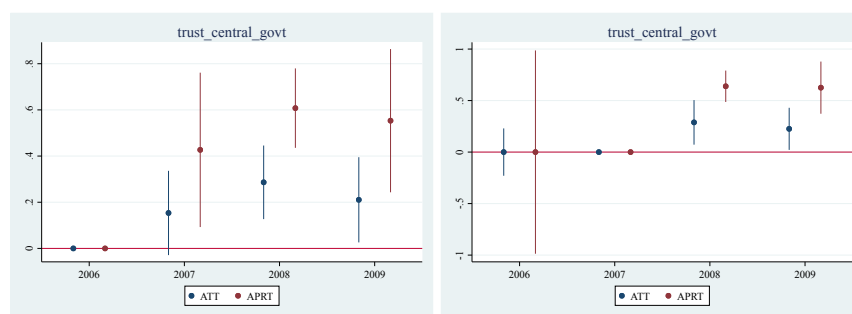
$$y_{icp} = \mu_{-1}^{(s)} + \sum_{j \neq -1} \alpha_j^{(s)} \mathbb{1}(c = s + j) + e_{icp}, \quad p \in P_s, \quad (33)$$

where P_s is the subset of provinces that introduced the new curriculum in year s . Then, as in (30), the estimator of $\theta(s, s + j)$, where $s \in \{2007, 2008, 2009\}$, is given by

$$\hat{\theta}(s, s + j) = \frac{\hat{\alpha}_j^{(s)} - \hat{\alpha}_j^{(2010)}}{1 - \hat{\mu}_{-1}^{(s)} - \hat{\alpha}_j^{(2010)'}}$$

where a hat denotes an estimator of the corresponding parameter. When $j \geq 0$, $\hat{\theta}(s, s + j)$ measures the average persuasion rate on the treated; whereas $j < -1$, it can be used for checking pre-treatment difference-in-differences.

FIGURE 1. Average Persuasion Rates on the Treated by Year of Treatment Adoption

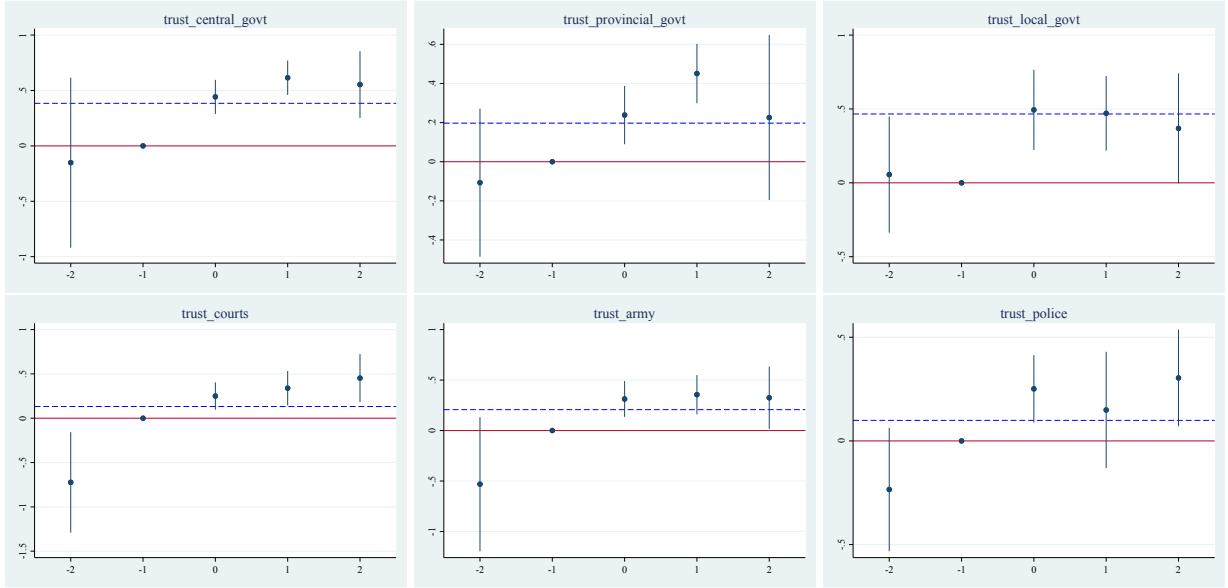


Notes: The left and right panels of the figure show estimates of both the average treatment effect on the treated (ATT) and the average persuasion rates on the treated (APRT) for $S = 2007$ and $S = 2008$ groups, respectively. The vertical lines represent 95% pointwise confidence intervals.

As an illustration, we consider the first outcome variable in [Cantoni et al. \(2017\)](#), the first row of Table 3), namely, “Trust: central government”. The left panel of Figure 1 shows the estimates of both ATT and APRT for the $S = 2007$ group. The X-axis shows different values of $s + j$, and by definition, 2006 corresponds to the null effect. The vertical lines represent 95% pointwise confidence intervals using the robust standard errors. They are obtained via the delta method, while combining the $S = s$ and $S = 2010$ regression models as a seemingly unrelated regression model with clustered dependence at the province level for each $S = s$ regression in (33). Not surprisingly, the APRT estimates are larger than the ATT estimates, as the former focuses on a more relevant subpopulation. It is remarkable to note that the 2007 ATT estimate is insignificant while the 2007 APRT estimate is significant, indicating that it is important to look at APRT as the key parameter. The right panel of Figure 1 shows the same estimates for the $S = 2008$ group. In this group, the 2006 estimates correspond to the pre-treatment period and are close to zero and statistically insignificant. Overall, both panels indicate that the impacts of the textbook reform were substantial on students’ trust on central government.

We do not report the estimates for the $S = 2009$ group separately as they turn out to be imprecisely estimated and include only one post-treatment period. Instead, we now report the estimator of the event-study persuasion rates (ESPR), i.e., $\theta_{ESPR}(j)$. If $j \geq 0$, it

FIGURE 2. Event-Study Persuasion Rates (ESPR)



Notes: Each of the figures shows estimates of the event-study persuasion rates (ESPR) for each horizon level. The vertical lines represent 95% pointwise confidence intervals.

can be obtained by

$$\hat{\theta}_{ESPR}(j) = \frac{\sum_{s=2007}^{2009-j} \hat{\mathbb{P}}(S = s) \{ \hat{\alpha}_j^{(s)} - \hat{\alpha}_j^{(2010)} \}}{\sum_{s=2007}^{2009-j} \hat{\mathbb{P}}(S = s) \{ 1 - \hat{\mu}_{-1}^{(s)} - \hat{\alpha}_j^{(2010)} \}}.$$

In addition, we also report the following estimate as a check for the pre-treatment period:

$$\hat{\theta}_{ESPR}(-2) = \frac{\sum_{s=2008}^{2009} \hat{\mathbb{P}}(S = s) \{ \hat{\alpha}_{-2}^{(s)} - \hat{\alpha}_{-2}^{(2010)} \}}{\sum_{s=2008}^{2009} \hat{\mathbb{P}}(S = s) \{ 1 - \hat{\mu}_{-1}^{(s)} - \hat{\alpha}_{-2}^{(2010)} \}}.$$

Figure 2 shows the ESPR estimates at each horizon level for the first 6 outcome variables from [Cantoni et al. \(2017\)](#), see Panel A in Table 3). The blue dashed line corresponds to the persuasion rate reported in [Cantoni et al. \(2017\)](#). It can be seen that our estimates are slightly larger, possibly and partly due to the fact that we allow for heterogeneous treatment effects, while carefully constructing control groups and dropping provinces that introduced the new curriculum before 2007. As before, pointwise 95% confidence intervals are plotted based on the robust standard errors, clustered at the province level for each $S = s$ regression model. On one hand, there is no significant effect in the pre-treatment period

except for the outcome “Trust: courts”. On the other hand, most of post-treatment period effects are significantly large. Overall, we observe that the main findings in [Cantoni et al. \(2017\)](#) are re-confirmed with our analysis, providing further evidence on their conclusion that “studying the new curriculum led to more positive views of China’s governance”.

10. CONCLUSIONS

We have developed an econometric framework for learning the average persuasion rates on the treated using difference-in-differences, providing detailed discussions on identification, estimation, and inference. While our framework is designed for panel data, it can be extended to repeated cross-sectional data. However, the specific modifications required for such an extension have not been fully explored in this paper. Investigating these extensions is a promising direction for future research.

APPENDIX A. PROOFS OF THE RESULTS IN THE MAIN TEXT.

Proof of Lemma 1: We generally have

$$\begin{aligned} & \mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\} - \mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\} \\ &= \mathbb{P}\{Y_1(0) = 0, Y_1(1) = 1 \mid D_1 = 1, X = x\} - \mathbb{P}\{Y_1(0) = 1, Y_1(1) = 0 \mid D_1 = 1, X = x\}. \end{aligned}$$

However, $\mathbb{P}\{Y_1(0) = 1, Y_1(1) = 0 \mid D_1 = 1, X = x\} = 0$ under assumption **B**. If assumption **B** is violated, then the claim follows from the lower bound of Fréchet-Hoeffding inequalities: see lemma 3 in Online Appendix [S-1](#) for more detail. \square

Proofs of Lemma 2: Suppose that assumption **B** holds for all $x \in \mathcal{X}$. Then, $\theta_c(x) = \theta_{cL}(x)$ and $\theta_c^{(r)}(x) = \theta_{cL}^{(r)}(x)$ for all $x \in \mathcal{X}$ by lemma 1. Therefore, multiplying the density of X at x given $D_1 = 1, Y_1(0) = 0$ to both sides of $\theta_c(x) = \theta_{cL}(x)$ and integrating with respect to x shows $\theta = \theta_L$. The other case of $\theta^{(r)} = \theta_L^{(r)}$ is similar. \square

Proof of Theorem 1: It follows from assumption **A** and equation (5) and the fact that $Y_1(d) = Y_1$ when we condition on $D_1 = d$. \square

Proof of Corollary 1: It follows from lemma 1 and theorem 1, and the Bayes rule. \square

Corollary 2 in Online Appendix S-1 contains extended results on sharp identified bounds when assumption B can be violated.

Proof of Theorem 2: Under assumption D, the coefficients in equation (8) are identified by OLS, where $D_{i1} = G_i$. Therefore, the denominator of $\bar{\theta}_L$ is

$$\begin{aligned} & \{1 - \mathbb{E}(Y_{i0} \mid D_{i1} = 1)\} - \{\mathbb{E}(Y_{i1} \mid D_{i1} = 0) - \mathbb{E}(Y_{i0} \mid D_{i1} = 0)\} \\ &= \{1 - \mathbb{E}(Y_{it} \mid G_i = 1, t = 0)\} - \{\mathbb{E}(Y_{it} \mid G_i = 0, t = 1) - \mathbb{E}(Y_{it} \mid G_i = 0, t = 0)\} \\ &= 1 - (\gamma_0 + \gamma_1) - (\gamma_0 + \gamma_2) + \gamma_0 = 1 - \gamma_0 - \gamma_1 - \gamma_2, \end{aligned}$$

and that of $\bar{\theta}_L^{(r)}$ is $\mathbb{E}(Y_{i1} \mid D_{i1} = 1) = \mathbb{E}(Y_{it} \mid G_i = 1, t = 1) = \gamma_0 + \gamma_1 + \gamma_2 + \gamma$. The numerators of θ_L and $\theta_L^{(r)}$ are similar. \square

Proof of Theorem 3: Let $\mathbb{E}(D_{i1}) = q$, and consider the denominator of the middle expression in equation (10): i.e.,

$$\frac{\text{Cov}(\tilde{Y}_{i1} - Y_{i0}, D_{i1})}{\mathbb{V}(D_{i1})} = \frac{\mathbb{E}(\tilde{Y}_{i1}D_{i1}) - \mathbb{E}(Y_{i0}D_{i1}) - q\{\mathbb{E}(\tilde{Y}_{i1}) - \mathbb{E}(Y_{i0})\}}{q(1-q)},$$

which is equal to

$$\begin{aligned} & \frac{\mathbb{E}(\tilde{Y}_{i1}D_{i1}) - \mathbb{E}(Y_{i0}D_{i1}) - q[\mathbb{E}(\tilde{Y}_{i1}D_{i1}) + \mathbb{E}\{\tilde{Y}_{i1}(1 - D_{i1})\} - \mathbb{E}(Y_{i0}D_{i1}) - \mathbb{E}\{Y_{i0}(1 - D_{i1})\}]}{q(1-q)} \\ &= \frac{(1-q)\{\mathbb{E}(\tilde{Y}_{i1}D_{i1}) - \mathbb{E}(Y_{i0}D_{i1})\} - q[\mathbb{E}\{\tilde{Y}_{i1}(1 - D_{i1})\} - \mathbb{E}\{Y_{i0}(1 - D_{i1})\}]}{q(1-q)} \\ &= \{\mathbb{E}(\tilde{Y}_{i1} \mid D_{i1} = 1) - \mathbb{E}(Y_{i0} \mid D_{i1} = 1)\} - \{\mathbb{E}(\tilde{Y}_{i1} \mid D_{i1} = 0) - \mathbb{E}(Y_{i0} \mid D_{i1} = 0)\} \\ &= 1 - \mathbb{E}(Y_{i0} \mid D_{i1} = 1) - \mathbb{E}(Y_{i1} \mid D_{i1} = 0) + \mathbb{E}(Y_{i0} \mid D_{i1} = 0), \end{aligned}$$

which is the denominator of $\bar{\theta}_L$. The numerator of $\bar{\theta}_L$ is similar. Finally, for the denominator of $\bar{\theta}_L^{(r)}$, just note that $\text{Cov}(Y_{i1}D_{i1}, D_{i1}) = \mathbb{E}(Y_{i1} \mid D_{i1} = 1)\mathbb{V}(D_{i1})$. \square

Proof of Theorem 4: Equivalence between the right-hand side of equation (17) and that of equation (18) is a simple algebraic result. The fact that $F_{DID}(Y_0, Y_1, D_1, X)$ is the semi-parametrically efficient influence function follows from lemma 9 in Online Appendix S-5

and the fact that $F_{num}(Y_0, Y_1, D_1, X)$ and $F_{den}(Y_0, Y_1, D_1, X)$ are in the tangent space \mathcal{T} described in lemma 5 in Online Appendix S-4: see the expressions of $F_{num}(Y_0, Y_1, D_1, X)$ and $F_{den}(Y_0, Y_1, D_1, X)$ given before the proofs of lemmas 6 and 7 in Online Appendix S-5. Finally, asymptotic normality in (19) follows from Theorem 2.1 of Newey (1994). Specifically, the scores are given in equation (42) in Online Appendix S-4, and they form a linear space. Also, any mean-zero function $s(Y_0, Y_1, D_1, X) = s_{000}(X)(1 - Y_0)(1 - Y_1)(1 - D_1) + \dots + s_{111}(X)Y_0Y_1D_1$ can be exactly matched with a score in the form of equation (42) in Online Appendix S-4. Therefore, an asymptotically linear estimator of $\bar{\theta}_L$ must have a unique influence function, which will coincide with F_{DID} by Theorem 2.1 of Newey (1994).

Proof of Theorem 5: It is useful to define a few immediate objects of interest. Let

$$\begin{aligned}\theta_{cL}(s, t \mid x) &:= \frac{\mathbb{P}\{Y_t(s) = 1 \mid S = s, X = x\} - \mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X = x\}}{1 - \mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X = x\}}, \\ \theta_L(s, t) &:= \frac{\mathbb{E}[\mathbb{P}\{Y_t(s) = 1 \mid S = s, X\} - \mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X\} \mid S = s]}{1 - \mathbb{E}[\mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X\} \mid S = s]}, \\ \theta_{ESPR,L}(j) &:= \frac{\sum_{s=1}^{T-j} \theta_L(s, s+j) \mathbb{P}(S = s) [1 - \mathbb{E}\{\tau_{stagger}(s, s+j \mid X) \mid S = s\}]}{\sum_{s=1}^{T-j} \mathbb{P}(S = s) [1 - \mathbb{E}\{\tau_{stagger}(s, s+j \mid x) \mid S = s\}]}.\end{aligned}$$

As in the proof of Lemma 1, we generally have

$$\begin{aligned}\mathbb{P}\{Y_t(s) = 1 \mid S = s, X = x\} - \mathbb{P}\{Y_t(\infty) = 1 \mid S = s, X = x\} \\ = \mathbb{P}\{Y_t(s) = 1, Y_t(\infty) = 0 \mid S = s, X = x\} - \mathbb{P}\{Y_t(s) = 0, Y_t(\infty) = 1 \mid S = s, X = x\}.\end{aligned}$$

However, under assumption **G**, $\mathbb{P}\{Y_t(s) = 0, Y_t(\infty) = 1 \mid S = s, X = x\} = 0$. Therefore, $\theta_c(s, t \mid x) = \theta_{cL}(s, t \mid x)$. Furthermore, we can verify that $\theta(s, t) = \theta_L(s, t)$ and $\theta_{ESPR}(j) = \theta_{ESPR,L}(j)$ using the fact that for generic random variables A, B , and an event E , we generally have $\mathbb{E}(A \mid B \in E) = \mathbb{E}\{\mathbb{E}(A \mid B) \mathbb{1}(B \in E)\} / \mathbb{P}(B \in E)$.

Now the first identification result in equation (24) follows immediately from equation (23), because $Y_t(\infty) = Y_t$ and $Y_{s-1}(\infty) = Y_{s-1}$ when $S = \infty$, and $\mathbb{P}\{Y_{s-1}(\infty) = 1 \mid S = s, X\} = \mathbb{P}\{Y_{s-1}(s) = 1 \mid S = s, X\}$ by assumption **F** that rules out anticipation. The second result in equation (25) follows from an aggregation of $\theta_{cL}(s, s+j \mid X)$ by using the conditional

distribution of X given $Y_{s+j}(\infty) = 0$ and $S = s$. Finally, to obtain the third result, note that $\theta_{ESPR,L}(j)$ is identified by

$$\theta_{ESPR,L}(j) = \frac{\sum_{s=1}^{T-j} \theta_L(s, s+j) \mathbb{P}(S = s) [1 - \mathbb{E}\{\Psi_{stagger}(s, s+j | X) | S = s\}]}{\sum_{s=1}^{T-j} \mathbb{P}(S = s) [1 - \mathbb{E}\{\Psi_{stagger}(s, s+j | x) | S = s\}]} \quad (34)$$

Combining (34) with (25) yields equation (26). \square

Proof of Theorem 6: We are conditioning on $S_i = s$ or $S_i = \infty$ in running this regression. However, since $\mathbb{E}(\cdot | S_i = r, S_i = s \text{ or } S_i = \infty) = \mathbb{E}(\cdot | S_i = r)$ for $r \in \{s, \infty\}$ in general, the claim follows by the same reasoning as in the proof of theorem 2. \square

REFERENCES

- ABADIE, ALBERTO (2005): "Semiparametric Difference-in-Differences Estimators," *Review of Economic Studies*, 72 (1), 1–19.
- ACKERBERG, D., X. CHEN, J. HAHN, AND Z. LIAO (2014): "Asymptotic efficiency of semiparametric two-step GMM," *Review of Economic Studies*, 81 (3), 919–943.
- ARKHANGELSKY, DMITRY, SUSAN ATHEY, DAVID A HIRSHBERG, GUIDO W IMBENS, AND STEFAN WAGER (2021): "Synthetic difference-in-differences," *American Economic Review*, 111 (12), 4088–4118.
- ATHEY, SUSAN AND GUIDO W IMBENS (2006): "Identification and inference in nonlinear difference-in-differences models," *Econometrica*, 74 (2), 431–497.
- (2022): "Design-based analysis in difference-in-differences settings with staggered adoption," *Journal of Econometrics*, 226 (1), 62–79.
- BAYES, THOMAS (1763): "LII. An essay towards solving a problem in the doctrine of chances. By the late Rev. Mr. Bayes, F.R.S. communicated by Mr. Price, in a letter to John Canton, A.M.F.R.S.," *Philosophical Transactions of the Royal Society*, 53, 370–418.
- BLANDHOL, CHRISTINE, JOHN BONNEY, MAGNE MOGSTAD, AND ALEXANDER TORGOVITSKY (2022): "When is TSLS Actually LATE?" Working Paper 29709, National Bureau of Economic Research.
- BURSZTYN, LEONARDO AND DAVID Y. YANG (2022): "Misperceptions About Others," *Annual Review of Economics*, 14, 425–452.

- CALLAWAY, BRANTLY AND PEDRO H.C. SANT'ANNA (2021): "Difference-in-Differences with Multiple Time Periods," *Journal of Econometrics*, 225 (2), 200–230.
- CANTONI, DAVIDE, YUYU CHEN, DAVID Y YANG, NOAM YUCHTMAN, AND Y JANE ZHANG (2017): "Curriculum and ideology," *Journal of Political Economy*, 125 (2), 338–392.
- CHERNOZHUKOV, VICTOR, JUAN CARLOS ESCANCIANO, HIDEHIKO ICHIMURA, WHITNEY K. NEWEY, AND JAMES M. ROBINS (2022): "Locally Robust Semiparametric Estimation," *Econometrica*, 90 (4), 1501–1535.
- DAWID, A. PHILIP, DAVID L. FAIGMAN, AND STEPHEN E. FIENBERG (2014): "Fitting Science Into Legal Contexts: Assessing Effects of Causes or Causes of Effects?" *Sociological Methods & Research*, 43 (3), 359–390.
- DAWID, A. PHILIP AND MONICA MUSIO (2022): "Effects of Causes and Causes of Effects," *Annual Review of Statistics and Its Application*, 9 (Volume 9, 2022), 261–287.
- DE CHAISEMARTIN, CLÉMENT AND XAVIER D'HAULTFOEUILLE (2020): "Two-way fixed effects estimators with heterogeneous treatment effects," *American Economic Review*, 110 (9), 2964–2996.
- DELLAVIGNA, STEFANO AND MATTHEW GENTZKOW (2010): "Persuasion: empirical evidence," *Annual Review of Economics*, 2, 643–669.
- DELLAVIGNA, STEFANO AND ETHAN KAPLAN (2007): "The Fox News effect: Media bias and voting," *Quarterly Journal of Economics*, 122 (3), 1187–1234.
- DE CHAISEMARTIN, CLÉMENT AND XAVIER D'HAULTFOEUILLE (2022): "Two-way fixed effects and differences-in-differences with heterogeneous treatment effects: a survey," *The Econometrics Journal*, 26 (3), C1–C30.
- DING, PENG AND FAN LI (2019): "A Bracketing Relationship between Difference-in-Differences and Lagged-Dependent-Variable Adjustment," *Political Analysis*, 27 (4), 605–615.
- FREYALDENHOVEN, SIMON, CHRISTIAN HANSEN, JORGE PÉREZ PÉREZ, AND JESSE M SHAPIRO (2021): "Visualization, Identification, and Estimation in the Linear Panel Event-Study Design," Working Paper 29170, National Bureau of Economic Research.
- FREYALDENHOVEN, SIMON, CHRISTIAN HANSEN, AND JESSE M. SHAPIRO (2019): "Pre-event Trends in the Panel Event-Study Design," *American Economic Review*, 109 (9), 3307–38.
- GOLDSMITH-PINKHAM, PAUL, PETER HULL, AND MICHAL KOLESÁR (2022): "Contamination Bias in Linear Regressions," Working Paper 30108, National Bureau of Economic Research.

- GOODMAN-BACON, ANDREW (2021): "Difference-in-differences with variation in treatment timing," *Journal of Econometrics*, 225 (2), 254–277.
- HAHN, JINYONG (1998): "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects," *Econometrica*, 66 (2), 315–331.
- HAINMUELLER, JENS (2012): "Entropy Balancing for Causal Effects: A Multivariate Reweighting Method to Produce Balanced Samples in Observational Studies," *Political Analysis*, 20 (1), 25–46.
- HECKMAN, JAMES J, JEFFREY SMITH, AND NANCY CLEMENTS (1997): "Making the most out of programme evaluations and social experiments: Accounting for heterogeneity in programme impacts," *Review of Economic Studies*, 64 (4), 487–535.
- HIRANO, KEISUKE, GUIDO W. IMBENS, AND GEERT RIDDER (2003): "Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score," *Econometrica*, 71 (4), 1161–1189.
- IMBENS, GUIDO W AND JOSHUA D. ANGRIST (1994): "Identification and estimation of local average treatment effects," *Econometrica*, 62 (2), 467–475.
- JI, WENLONG, LIHUA LEI, AND ASHER SPECTOR (2023): "Model-Agnostic Covariate-Assisted Inference on Partially Identified Causal Effects," ArXiv:2310.08115 [econ.EM], available at <https://arxiv.org/abs/2310.08115>.
- JUN, SUNG JAE AND SOKBAE LEE (2023): "Identifying the effect of persuasion," *Journal of Political Economy*, 131 (8), 2032–2058.
- KAJI, TETSUYA AND JIANFEI CAO (2023): "Assessing Heterogeneity of Treatment Effects," ArXiv:2306.15048 [econ.EM], available at <https://arxiv.org/abs/2306.15048>.
- LADD, JONATHAN MCDONALD AND GABRIEL S. LENZ (2009): "Exploiting a Rare Communication Shift to Document the Persuasive Power of the News Media," *American Journal of Political Science*, 53 (2), 394–410.
- MANSKI, CHARLES F. (1997): "Monotone treatment response," *Econometrica*, 1311–1334.
- NEWAY, WHITNEY K (1994): "The asymptotic variance of semiparametric estimators," *Econometrica*, 1349–1382.
- PEARL, JUDEA (1999): "Probabilities of causation: three counterfactual interpretations and their identification," *Synthese*, 121 (1-2), 93–149.
- POSSEBOM, VITOR AND FLAVIO RIVA (2024): "Probability of Causation with Sample Selection: A Reanalysis of the Impacts of Jóvenes en Acción on Formality," *Journal of Business & Economic*

Statistics, in press, eprint available at <https://arxiv.org/abs/2210.01938>.

ROTH, JONATHAN AND PEDRO HC SANT'ANNA (2023): "When is parallel trends sensitive to functional form?" *Econometrica*, 91 (2), 737–747.

ROTH, JONATHAN, PEDRO HC SANT'ANNA, ALYSSA BILINSKI, AND JOHN POE (2023): "What's trending in difference-in-differences? A synthesis of the recent econometrics literature," *Journal of Econometrics*, 235 (2), 2218–2244.

SANT'ANNA, PEDRO HC AND JUN ZHAO (2020): "Doubly robust difference-in-differences estimators," *Journal of Econometrics*, 219 (1), 101–122.

SUN, LIYANG AND SARAH ABRAHAM (2021): "Estimating dynamic treatment effects in event studies with heterogeneous treatment effects," *Journal of Econometrics*, 225 (2), 175–199.

SUN, LIYANG AND JESSE M. SHAPIRO (2022): "A Linear Panel Model with Heterogeneous Coefficients and Variation in Exposure," *Journal of Economic Perspectives*, 36 (4), 193–204.

WOOLDRIDGE, JEFFREY M (2023): "Simple approaches to nonlinear difference-in-differences with panel data," *The Econometrics Journal*, 26 (3), C31–C66.

XU, YIQING, ANQI ZHAO, AND PENG DING (2024): "Factorial Difference-in-Differences," ArXiv:2407.11937 [stat.ME], available at <https://arxiv.org/abs/2407.11937>.

YAMAMOTO, TEPPEI (2012): "Understanding the Past: Statistical Analysis of Causal Attribution," *American Journal of Political Science*, 56 (1), 237–256.

YU, ARTHUR ZEYANG (2023): "A Binary IV Model for Persuasion: Profiling Persuasion Types among Compliers," Working Paper, available at https://arthurzeyangyu.github.io/jmp/yu_2023local.pdf.

ZHANG, CHAO, ZHI GENG, WEI LI, AND PENG DING (2024): "Identifying and bounding the probability of necessity for causes of effects with ordinal outcomes," ArXiv:2411.01234 [math.ST], available at <https://arxiv.org/abs/2411.01234>.

Online Appendices to “Learning the Effect of Persuasion via Difference-in-Differences”

Sung Jae Jun

Sokbae Lee

Penn State University

Columbia University

APPENDIX S-1. THE CASE OF THE BACKLASH

Without making assumption **B**, the Fréchet–Hoeffding inequality yields bounds on $\theta_c(x)$ and $\theta_c^{(r)}(x)$ in terms of the marginal probabilities of the potential outcomes conditional on $X = x$. Further, $\theta_c(x)$ and $\theta_c^{(r)}(x)$ are linearly dependent in that for all $x \in \mathcal{X}$, we have

$$\theta_c^{(r)}(x) = \frac{\{1 - \tau_c(x)\}\theta_c(x)}{\mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\}}. \quad (35)$$

Therefore, assuming that $\tau_c(x)$ is known, it suffices to have bounds on $\theta_c(x)$ to have bounds on $\theta_c^{(r)}(x)$.

Below we first derive bounds on $\theta_c(x)$ without using assumption **B** that depend only on the marginal probabilities of the potential outcomes given $D_1 = 1$ and $X = x$. Then, bounds on $\theta_c^{(r)}(x)$ will follow from the bounds on $\theta_c(x)$ and equation (35), or vice versa.

Define

$$\begin{aligned} \mathcal{B}(x) &:= \left\{ (p, q) \in \mathbb{R}^2 : \max\{0, \theta_{cL}(x)\} \leq p \leq \min\{\theta_{cU}(x), 1\}, q = \alpha(x)p \right\} \\ &= \left\{ (p, q) \in \mathbb{R}^2 : p = q/\alpha(x), \max\{0, \theta_{cL}^{(r)}(x)\} \leq q \leq \min\{\theta_{cU}^{(r)}(x), 1\} \right\}, \end{aligned}$$

where $\alpha(x) := \{1 - \tau_c(x)\}/\mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\}$, and

$$\theta_{cU}(x) := \frac{\mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\}}{1 - \tau_c(x)}, \quad \theta_{cU}^{(r)}(x) := \frac{1 - \tau_c(x)}{\mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\}}.$$

We then have the following lemma.

Lemma 3. Suppose that assumption **A** holds. For all $x \in \mathcal{X}$, we have $(\theta_c(x), \theta_c^{(r)}(x)) \in \mathcal{B}$, which is sharp based on the information of $\mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\}$ and $\mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\}$.

Proof. It follows from equation (35) and an application of Lemma I1 in Appendix I of **Jun and Lee (2023)** by setting $\mathbb{P}^*(A \cap B)$ in the lemma to be $\mathbb{P}\{Y_1(0) \in A, Y_1(1) \in B \mid D_1 = 1, X = x\}$. \square

The set \mathcal{B} is a line that describes the bounds on $(\theta_c(x), \theta_c^{(r)}(x))$ we can obtain from the marginal probabilities of the potential outcomes (given $D_1 = 1, X = x$) without using assumption **B**. Therefore, lemmas **1** and **3** show that rescaling CATT as in the end point $(\theta_{cL}(x), \theta_{cL}^{(r)}(x))$ provides useful causal parameters to consider. Conditioning on $X = x$, $(\theta_{cL}(x), \theta_{cL}^{(r)}(x))$ exactly corresponds to the persuasion rates on the treated that we are interested in if there is no backlash, while they serve as conservative measures in general, even if the backlash effect is a concern.

Aggregating X with appropriate conditioning shows that θ_L and $\theta_L^{(r)}$ are valid lower bounds on θ and $\theta^{(r)}$, respectively: i.e., even if there is a concern about the backlash so that assumption **B** may be violated, θ_L and $\theta_L^{(r)}$ still offer conservative measures of θ and $\theta^{(r)}$, respectively, although it may not be sharp in general. In order to be more precise on this issue, let

$$\mathbb{L} := \mathbb{E}\{\mathbb{1}_{\mathcal{X}_L}(X)\text{CATT}(X) \mid D_1 = 1\},$$

$$\mathbb{U} := \mathbb{E}[\mathbb{1}_{\mathcal{X}_U}(X)\mathbb{P}(Y_1 = 1 \mid D_1 = 1, X) + \{1 - \mathbb{1}_{\mathcal{X}_U}(X)\}\{1 - \tau_c(X)\} \mid D_1 = 1],$$

where \mathcal{X}_L and \mathcal{X}_U are defined by

$$\mathcal{X}_L := \{x \in \mathcal{X} : \mathbb{P}(Y_1 = 1 \mid D_1 = 1, X = x) - \tau_c(x) \geq 0\},$$

$$\mathcal{X}_U := \{x \in \mathcal{X} : \tau_c(x) + \mathbb{P}(Y_1 = 1 \mid D_1 = 1, X = x) \leq 1\}.$$

Define

$$\mathcal{B} := \left\{ (p, q) \in \mathbb{R}^2 : \theta_L^* \leq p \leq \theta_U^*, q = \alpha p \right\} = \left\{ (p, q) \in \mathbb{R}^2 : p = q/\alpha, \theta_L^{(r)*} \leq q \leq \theta_U^{(r)*} \right\},$$

where $\alpha = \mathbb{E}\{1 - \tau_c(X) \mid D_1 = 1\} / \mathbb{P}(Y_1 = 1 \mid D_1 = 1)$, and

$$\theta_L^* := \mathbb{L} / \mathbb{E}\{1 - \tau_c(X) \mid D_1 = 1\} \quad \text{and} \quad \theta_U^* := \mathbb{U} / \mathbb{E}\{1 - \tau_c(X) \mid D_1 = 1\},$$

$$\theta_L^{(r)*} := \mathbb{L} / \mathbb{P}(Y_1 = 1 \mid D_1 = 1) \quad \text{and} \quad \theta_U^{(r)*} := \mathbb{U} / \mathbb{P}(Y_1 = 1 \mid D_1 = 1).$$

Lemma 4. *Suppose that assumption A holds for all $x \in \mathcal{X}$. We then have $(\theta, \theta^{(r)}) \in \mathcal{B}$, which is sharp based on the information of $\mathbb{P}\{Y_1(1) = 1 \mid D_1 = 1, X = x\}$ and $\mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\}$ for all $x \in \mathcal{X}$, and the distribution of X given $D_1 = 1$.*

Proof. We will only show the bounds on θ : the set \mathcal{B} will follow from them and the linear relationship between θ and $\theta^{(r)}$, i.e., $\theta^{(r)} = \theta \mathbb{E}\{1 - \tau_c(X) \mid D_1 = 1\} / \mathbb{P}(Y_1 = 1 \mid D_1 = 1)$.

First, we note that

$$\mathcal{X}_L = \{x \in \mathcal{X} : \theta_{cL}(x) \geq 0\} \quad \text{and} \quad \mathcal{X}_U = \{x \in \mathcal{X} : \theta_{cU}(x) \leq 1\}$$

by definition. Therefore, we can equivalently write the sharp bounds in lemma 3 as

$$\mathbb{1}_{\mathcal{X}_L}(x) \theta_{cL}(x) \leq \theta_c(x) \leq \mathbb{1}_{\mathcal{X}_U}(x) \theta_{cU}(x) + 1 - \mathbb{1}_{\mathcal{X}_U}(x). \quad (36)$$

We then use the fact

$$f\{x \mid Y_1(0) = 0, D_1 = 1\} = \frac{\mathbb{P}\{Y_1(0) = 0 \mid D_1 = 1, X = x\} f(x \mid D_1 = 1)}{\mathbb{P}\{Y_1(0) = 0 \mid D_1 = 1\}} \quad (37)$$

by Bayes' rule. Specifically, multiplying the conditional density in equation (37) to both sides of the inequalities in equation (36) and integrating shows the claim. \square

Here, $\tau_c(\cdot)$ is the only unidentified object so that θ_L^* and θ_U^* will be identified if $\tau_c(\cdot)$ is identified. The set \mathcal{X}_L represents the values x of X such that $\text{CATT}(x) \geq 0$: i.e., it is the set of x such that the lower bound on $\theta_c(x)$ in lemma 3 is nontrivial. Similarly, \mathcal{X}_U is the set of the values $x \in \mathcal{X}$ such that the upper bound on $\theta_c(x)$ in lemma 3 is not trivial: or,

equivalently, it is the set of values $x \in \mathcal{X}$ such that the upper bound on $\theta_c^{(r)}(x)$ is trivial. Indeed, the condition that defines \mathcal{X}_U can be equivalently expressed as

$$\mathbb{P}\{Y_1(0) = 1, Y_1(1) = 1 \mid D_1 = 1, X = x\} \leq \mathbb{P}\{Y_1(0) = 0, Y_1(1) = 0 \mid D_1 = 1, X = x\}.$$

If this inequality is not satisfied so that there are too many ‘voters’ who are characterized by $X = x$ and who would vote for the party the media publicly endorses no matter what, then the upper bound on $\theta_c(x)$ based on the marginals of the potential outcomes will be just trivial, i.e., $\theta_{cU}(x) = 1$.

If assumption **B** holds, then we have $\text{CATT}(X) \geq 0$ almost surely, and therefore $\theta_L^* = \theta_L$ as well as $\theta_L^{(r)*} = \theta_L^{(r)}$ will follow. Therefore, lemmas 2 and 4 show that θ_L^* and $\theta^{(r)*}$ will be interesting parameters to consider: they are sharp lower bounds on θ and $\theta^{(r)}$, respectively, in general, while they are exactly equal to θ and $\theta^{(r)}$ under monotonicity.

However, $(\theta_L^*, \theta_L^{(r)*})$ is a more difficult parameter than $(\theta_L, \theta_L^{(r)})$ because the former contains $\mathbb{1}_{\mathcal{X}_L}(X)$ that depends on unknown objects in a nonsmooth way. In contrast, $(\theta_L, \theta_L^{(r)})$ can be estimated in a more straightforward manner, as long as $\tau_c(\cdot)$ is identified. Therefore, we consider θ_L and $\theta_L^{(r)}$ the aggregate parameters of interest. Since $\max(0, \theta_L) \leq \theta_L^*$ in general, θ_L is always a robust lower bound on θ ; the same comment applies to $\theta_L^{(r)}$ as well. If $\text{CATT}(x)$ is nonnegative for almost all $x \in \mathcal{X}$, then θ_L and $\theta_L^{(r)}$ are the sharp lower bounds θ_L^* and $\theta_L^{(r)*}$, respectively. Further, if assumption **B** holds, then $\theta_L = \theta_L^* = \theta$, and $\theta_L^{(r)} = \theta_L^{(r)*} = \theta^{(r)}$.

Partial identification under assumption **C** without using assumption **B** is largely uneventful. For example, lemma 4 and theorem 1 show the joint sharp identified set of CPRT and R-CPRT when assumption **B** is violated. For the aggregated parameters, the joint sharp identified set requires aggregation over an unknown subset of the support of X in general. Below we clarify this issue, and we make a formal statement about identification of the aggregated parameters.

Define

$$\begin{aligned}\bar{\theta}_L^* &:= \frac{\bar{\mathbb{L}}}{\mathbb{E}\{1 - \Psi(X) \mid D_1 = 1\}}, & \bar{\theta}_U^* &:= \frac{\bar{\mathbb{U}}}{\mathbb{E}\{1 - \Psi(X) \mid D_1 = 1\}}, \\ \bar{\theta}_L^{(r)*} &:= \frac{\bar{\mathbb{L}}}{\mathbb{E}\{\Pi_1(1, X) \mid D_1 = 1\}}, & \bar{\theta}_U^{(r)*} &:= \frac{\bar{\mathbb{U}}}{\mathbb{E}\{\Pi_1(1, X) \mid D_1 = 1\}},\end{aligned}$$

where

$$\begin{aligned}\bar{\mathbb{L}} &:= \mathbb{E}[\mathbb{1}_{\mathcal{X}_L^*}(X) \{\Pi_1(1, X) - \Psi(X)\} \mid D_1 = 1], \\ \bar{\mathbb{U}} &:= \mathbb{E}[\mathbb{1}_{\mathcal{X}_U^*}(X) \Pi_1(1, X) + \{1 - \mathbb{1}_{\mathcal{X}_U^*}(X)\} \{1 - \Psi(X)\} \mid D_1 = 1],\end{aligned}$$

with

$$\mathcal{X}_L^* := \{x \in \mathcal{X} : \Psi(x) \leq \Pi_1(1, x)\} \quad \text{and} \quad \mathcal{X}_U^* := \{x \in \mathcal{X} : \Psi(x) \leq 1 - \Pi_1(1, x)\}.$$

Here, $\bar{\theta}_L$, $\bar{\theta}_L^{(r)}$, $\bar{\theta}_L^*$, $\bar{\theta}_L^{(r)*}$, $\bar{\theta}_U^*$, and $\bar{\theta}_U^{(r)}$ are all directly identified from the data.

Corollary 2. *Suppose that assumptions **A** and **C** are satisfied for all $x \in \mathcal{X}$. Then, $(\theta_L, \theta_L^*, \theta_U^*)$ and $(\theta_L^{(r)}, \theta_L^{(r)*}, \theta_U^{(r)*})$ are identified by $(\bar{\theta}_L, \bar{\theta}_L^*, \bar{\theta}_U^*)$ and $(\bar{\theta}_L^{(r)}, \bar{\theta}_L^{(r)*}, \bar{\theta}_U^{(r)*})$, respectively. Therefore,*

- (i) *if assumptions **A** to **C** hold for all $x \in \mathcal{X}$, then $\theta = \theta_L$ and $\theta^{(r)} = \theta_L^{(r)}$ are point-identified by $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$, respectively;*
- (ii) *if assumptions **A** and **C** hold for all $x \in \mathcal{X}$, then the joint sharp identifiable set of $(\theta, \theta^{(r)})$ is given by the line connecting $[\bar{\theta}_L^*, \bar{\theta}_L^{(r)*}]$, and $[\bar{\theta}_U^*, \bar{\theta}_U^{(r)*}]$.*

Proof. Noting that $\mathcal{X}_L^* = \mathcal{X}_L$ and $\mathcal{X}_U^* = \mathcal{X}_U$ under assumptions **A** and **C** by theorem 1, the claim immediately follows from lemma 4 and theorem 1. \square

Here, $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$ are natural estimands to focus on. They are conservative measures of APRT and R-APRT in general, while they are exactly equal to APRT and R-APRT under monotonicity. Also, they are easier parameters to estimate than $\bar{\theta}_L^*$ or $\bar{\theta}_L^{(r)*}$, which depends on unknown objects in a nonsmooth way.

APPENDIX S-2. USING A KNOWN LINK FUNCTION

As we commented in section 3, assumption C does not allow for popular parametric models such as logit or probit. However, we can modify assumption C to introduce a link function, as long as the link function is pre-specified. Below is a modification of assumption C.

Assumption I (Parallel Trends). *For some known link function Λ on $[0, 1]$ that is strictly increasing and differentiable, and for all $x \in \mathcal{X}$, $\Lambda[\mathbb{P}\{Y_t(0) = 1 \mid D_1 = d, X = x\}]$ is separable into the sum of a time component and a treatment component: i.e., there exist functions G (of t, x) and H (of d, x) such that*

$$\mathbb{P}\{Y_t(0) = 1 \mid D_1 = d, X = x\} = \Lambda^{-1}\{G(t, x) + H(d, x)\}.$$

Differentiability of Λ will be useful for calculating the efficient influence function, not strictly necessary for identification. The choice of the link function Λ is a specification issue for the researcher: e.g., $\Lambda(s) = s$ is an obvious choice that we used throughout the main text. More generally, assumption I allows for the class of generalized linear models. For example, the logistic model with $\Lambda^{-1}(s) = \exp(s) / \{1 + \exp(s)\}$ and

$$\mathbb{P}\{Y_t(0) = 1 \mid D_1 = d, X = x\} = \Lambda^{-1}(\beta_0 + \beta_1 t + \beta_2 d + \beta_3 x + \beta_4^\top t x + \beta_5^\top d x)$$

does not satisfy assumption C, but it does satisfy assumption I.

In addition to the linear or logistic choice of Λ , it is worth considering $\Lambda^{-1}(s) = 1 - \exp(-s)$ with $s \geq 0$, i.e., the distribution function of the standard exponential distribution. This choice of the link function is for the case where the time and treatment component are multiplicatively separable, and therefore there is common growth as in [Wooldridge \(2023\)](#). Specifically, suppose that $\mathbb{P}\{Y_t(0) = 0 \mid D_1 = d, X = x\} = \tilde{G}(t, x)\tilde{H}(d, x)$ so that for all $x \in \mathcal{X}$,

$$\frac{\mathbb{P}\{Y_1(0) = 0 \mid D_1 = 1, X = x\}}{\mathbb{P}\{Y_0(0) = 0 \mid D_1 = 1, X = x\}} = \frac{\mathbb{P}\{Y_1(0) = 0 \mid D_1 = 0, X = x\}}{\mathbb{P}\{Y_0(0) = 0 \mid D_1 = 0, X = x\}}.$$

In this case, the choice of $\Lambda^{-1}(s) = 1 - \exp(-s)$ leads to

$$\Lambda[\mathbb{P}\{Y_t(0) = 1 \mid D_1 = d, X = x\}] = -\log \tilde{G}(t, x) - \log \tilde{H}(d, x).$$

Under assumption **I**, we have parallel trends with the transformation Λ : i.e.,

$$\begin{aligned} & \Lambda[\mathbb{P}\{Y_1(0) = 1 \mid D_1 = 1, X = x\}] - \Lambda[\mathbb{P}\{Y_0(0) = 1 \mid D_1 = 1, X = x\}] \\ &= \Lambda[\mathbb{P}\{Y_1(0) = 1 \mid D_1 = 0, X = x\}] - \Lambda[\mathbb{P}\{Y_0(0) = 1 \mid D_1 = 0, X = x\}]. \end{aligned} \quad (38)$$

Therefore, theorem **1** and corollaries **1** and **2** continue to hold with the modification of

$$\Psi(X) := \Lambda^{-1}[\Lambda\{\Pi_0(1, X)\} + \Lambda\{\Pi_1(0, X)\} - \Lambda\{\Pi_0(0, X)\}]. \quad (39)$$

Indeed, most of our results in the main text can be extended by using assumption **I** instead of assumption **C** with some exceptions: e.g., the plug-in approach that uses Ψ is valid for any choice of Λ , while the propensity-odds-weighting approach in equation (14) relies on the specific choice of $\Lambda(s) = s$.

APPENDIX S-3. FURTHER DISCUSSION ON CONTROLLING FOR THE PRE-TREATMENT OUTCOME AND UNCONFOUNDEDNESS

In this part of the appendix, we return to the issues discussed in section 3.2 and provide a more detailed discussion. Specifically, the unconfoundedness condition is an alternative assumption that has been used in the literature especially with cross-sectional data. When the pre-treatment outcome variable is available, it is natural to include it in the conditioning variables to define the unconfoundedness assumption. Below we first compare the unconfoundedness assumption with the parallel trend assumption. Let $Z := [Y_0, X^\top]^\top$, and consider the following assumptions.

Assumption J. *At time $t = 0$, no one is treated. At time $t = 1$, there is a constant $\epsilon > 0$ such that $\epsilon \leq \min[\mathbb{P}\{Y_1(0) = 0, D_1 = 1 \mid Z\}, \mathbb{P}\{Y_1(1) = 1, D_1 = 1 \mid Z\}]$ and $\mathbb{P}(D_1 = 1 \mid Z) \leq 1 - \epsilon$ with probability one.*

Assumption K. $Y_1(0)$ is independent of D_1 conditional on Z .

Assumption J is simply setting up the same environment as in the main text, but with Z in lieu of X (see assumption A). Assumption K imposes unconfoundedness given Z . We remark that assumption C holds with Z in lieu of X if and only if assumption K holds: this can be verified by using the fact that Z contains Y_0 .

Under assumptions J and K, both θ_L and $\theta_L^{(r)}$ are well-defined, and they are identified by

$$\begin{aligned}\tilde{\theta}_{L,Z} &:= \frac{\mathbb{E}\{D_1(Y_1 - Y_0)\} - \mathbb{E}\{D_1\mathbb{E}(Y_1 - Y_0|D_1 = 0, Z)\}}{\mathbb{E}\{D_1(1 - Y_0)\} - \mathbb{E}\{D_1\mathbb{E}(Y_1 - Y_0|D_1 = 0, Z)\}}, \\ \tilde{\theta}_{L,Z}^{(r)} &:= \frac{\mathbb{E}\{D_1(Y_1 - Y_0)\} - \mathbb{E}\{D_1\mathbb{E}(Y_1 - Y_0|D_1 = 0, Z)\}}{\mathbb{E}(D_1Y_1)},\end{aligned}\tag{40}$$

where we use the fact that Y_0 is included in Z . The expressions in equation (40) are reminiscent of $\bar{\theta}_L$ and $\bar{\theta}_L^{(r)}$ that identify θ_L and $\theta_L^{(r)}$ under assumption C: see corollary 1. Indeed,

$$\begin{aligned}\bar{\theta}_L &= \frac{\mathbb{E}\{D_1(Y_1 - Y_0)\} - \mathbb{E}\{D_1\mathbb{E}(Y_1 - Y_0|D_1 = 0, X)\}}{\mathbb{E}\{D_1(1 - Y_0)\} - \mathbb{E}\{D_1\mathbb{E}(Y_1 - Y_0|D_1 = 0, X)\}}, \\ \bar{\theta}_L^{(r)} &= \frac{\mathbb{E}\{D_1(Y_1 - Y_0)\} - \mathbb{E}\{D_1\mathbb{E}(Y_1 - Y_0|D_1 = 0, X)\}}{\mathbb{E}(D_1Y_1)}\end{aligned}\tag{41}$$

have the same forms as $\tilde{\theta}_{L,Z}$ and $\tilde{\theta}_{L,Z}^{(r)}$ except that they do not include Y_0 in the controls. Put differently, adding Y_0 to X and using the DID formula is an implementation of identifying ATT via assumption K instead of assumption C.

One might prefer the unconfoundedness condition after controlling for both X and Y_0 if it is important to avoid the functional form restriction of the parallel trend assumption. Alternatively, the parallel trend assumption might be favored since it would not require fully conditioning on Y_0 . Generally speaking, the required X under the parallel trend assumption could be different from X under the unconfoundedness assumption. If both are the same, which identification assumption to adopt is just reduced to the matter of whether to include Y_0 in the covariates or not.

Further, there is a simple testable condition under which it becomes moot to distinguish the two identification strategies. If Y_0 is independent of D_1 conditional on X , then it can

be shown that $\tilde{\theta}_{L,Z} = \bar{\theta}_L$ and $\tilde{\theta}_{L,Z}^{(r)} = \bar{\theta}_L^{(r)}$. Below we provide a more detailed discussion on this point. We first consider the following assumptions.

Assumption L. $Y_1(0)$ is independent of D_1 conditional on X .

Assumption M. Y_0 is independent of D_1 conditional on X .

Unlike Assumption **K**, Assumption **L** imposes unconfoundedness given X only. Assumptions **K** and **L** are not testable, but Assumption **M** is: Y_0 is observed for all observational units, but $Y_1(0)$ is not.

Let $\tilde{\theta}_{L,X}$ and $\tilde{\theta}_{L,X}^{(r)}$ be estimands that identify θ_L and $\theta_L^{(r)}$, respectively, under unconfoundedness given X (i.e., assumption **L**). That is, define

$$\tilde{\theta}_{L,X} := \frac{\mathbb{E}(D_1 Y_1) - \mathbb{E}\{D_1 \mathbb{E}(Y_1 | D_1 = 0, X)\}}{\mathbb{E}(D_1) - \mathbb{E}\{D_1 \mathbb{E}(Y_1 | D_1 = 0, X)\}}, \quad \tilde{\theta}_{L,X}^{(r)} := \frac{\mathbb{E}(D_1 Y_1) - \mathbb{E}\{D_1 \mathbb{E}(Y_1 | D_1 = 0, X)\}}{\mathbb{E}(D_1 Y_1)}.$$

Now, in view of (40), note that $\tilde{\theta}_{L,Z}$ and $\tilde{\theta}_{L,Z}^{(r)}$ can be equivalently written as $\tilde{\mathcal{N}} / [\tilde{\mathcal{N}} + \mathbb{E}\{D_1(1 - Y_1)\}]$ and $\tilde{\mathcal{N}} / \mathbb{E}(D_1 Y_1)$, respectively, where

$$\tilde{\mathcal{N}} := \mathbb{E} \left\{ D_1(Y_1 - Y_0) - (1 - D_1)(Y_1 - Y_0) \frac{\mathbb{P}(D_1 = 1 | Z)}{\mathbb{P}(D_1 = 0 | Z)} \right\}.$$

Recall that $\bar{\theta}_L = \mathcal{N} / [\mathcal{N} + \mathbb{E}\{D_1(1 - Y_1)\}]$ and $\bar{\theta}_L^{(r)} = \mathcal{N} / \mathbb{E}(D_1 Y_1)$, where \mathcal{N} is given in (15). Therefore, comparing $\tilde{\mathcal{N}}$ with \mathcal{N} shows that assumption **M** implies that $\bar{\theta}_L = \tilde{\theta}_{L,Z}$, and $\bar{\theta}_L^{(r)} = \tilde{\theta}_{L,Z}^{(r)}$. Also, it follows from

$$\mathbb{E}\{D_1 \mathbb{E}(Y_1 | D_1 = 0, Z)\} = \mathbb{E} \left\{ Y_1(1 - D_1) \frac{\mathbb{P}(D_1 = 1 | Z)}{\mathbb{P}(D_1 = 0 | Z)} \right\}$$

that assumption **M** is again sufficient to ensure that $\tilde{\theta}_{L,Z} = \tilde{\theta}_{L,X}$ as well as $\tilde{\theta}_{L,Z}^{(r)} = \tilde{\theta}_{L,X}^{(r)}$.

Our discussion so far can be summarized as in the following remark.

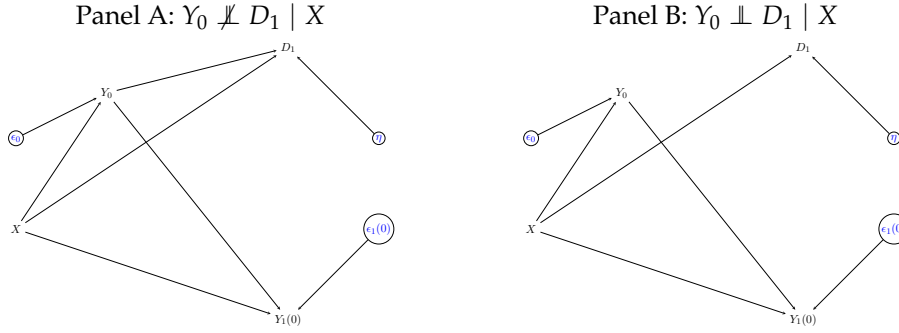
Remark 1. Suppose that assumption **J** holds.

- (i) If assumption **M** holds, then assumption **C** is the same as assumption **L**, while we have $\bar{\theta}_L = \tilde{\theta}_{L,Z} = \tilde{\theta}_{L,X}$ and $\bar{\theta}_L^{(r)} = \tilde{\theta}_{L,Z}^{(r)} = \tilde{\theta}_{L,X}^{(r)}$. Hence, θ_L and $\theta_L^{(r)}$ are identified by

$\bar{\theta}_L = \tilde{\theta}_{L,Z} = \tilde{\theta}_{L,X}$ and $\bar{\theta}_L^{(r)} = \tilde{\theta}_{L,Z}^{(r)} = \tilde{\theta}_{L,X}^{(r)}$, respectively, under either assumption **C** or assumption **K**.¹⁴

- (ii) If assumption **M** does not hold, then the researcher needs to take a stance among assumptions **C**, **K**, or **L**.

FIGURE 3. Examples of Data Generating Processes with $Z = [Y_0, X^T]^T$



Notes: Variables in circles are unobserved and they are independent. Assumption **K** holds in both cases.

Figure 3 illustrates some potential data generating processes. In Panel A, Y_0 is not independent of D_1 given X . Therefore, the researcher needs to take a stance. In this diagram, assumption **K** is satisfied, and therefore, we can identify $(\theta_L, \theta_L^{(r)})$ by $(\tilde{\theta}_{L,Z}, \tilde{\theta}_{L,Z}^{(r)})$. However, conditioning only on X is not sufficient to deliver independence of D_1 and $Y_1(0)$, so $\tilde{\theta}_{L,X}$ and $\tilde{\theta}_{L,X}^{(r)}$ are not valid estimands to identify θ_L and $\theta_L^{(r)}$. It is not clear from this diagram whether assumption **C** is satisfied or not though. In Panel B, D_1 is independent of Y_0 given X . Hence, it does not matter what stance the researcher takes: $\bar{\theta}_L$, $\tilde{\theta}_{L,Z}$, and $\tilde{\theta}_{L,X}$ are all equal to θ_L , while $\bar{\theta}_L^{(r)}$, $\tilde{\theta}_{L,Z}^{(r)}$, and $\tilde{\theta}_{L,X}^{(r)}$ are all equal to $\theta_L^{(r)}$. In the diagram, assumption **K** holds, and therefore, both assumptions **C** and **L** are satisfied as well.

The main takeaway from this discussion is that if Y_0 is independent of D_1 given X , then it is largely an unimportant question whether to control for Y_0 or not, or whether to rely on unconfoundedness or parallel trends. However, if Y_0 is not independent of D_1 given X , then the researcher needs to decide carefully which estimand to rely on to learn about θ_L or $\theta_L^{(r)}$.

¹⁴In fact, it is easy to verify that assumption **K** implies assumption **C** in this case.

It seems feasible to establish a bracketing relationship between the two identification approaches, drawing on the findings of [Ding and Li \(2019\)](#) regarding ATT; however, we leave this extension for future research.

APPENDIX S-4. DERIVATION OF THE TANGENT SPACE

We observe $(Y_0, Y_1, D_1, X^\top)^\top$: we also observe D_0 , but it is irrelevant for our discussion here, because $D_0 = 0$ with probability one by the setup. Let f be the density of X , and let $P(X) := \mathbb{P}(D_1 = 1 \mid X)$. For $d \in \{0, 1\}$, let $Q_d(X) := \mathbb{P}(Y_0 = 1 \mid D_1 = d, X)$. Further, for $d, y \in \{0, 1\}$, let $R_{dy}(X) := \mathbb{P}(Y_1 = 1 \mid D_1 = d, Y_0 = y, X)$. Then, the likelihood is the product of the following terms, while each line corresponds to one term:

$$\begin{aligned}
& f(X) \\
& P(X)^{D_1} \{1 - P(X)\}^{1-D_1} \\
& [Q_1(X)^{Y_0} \{1 - Q_1(X)\}^{1-Y_0}]^{D_1} \\
& [Q_0(X)^{Y_0} \{1 - Q_0(X)\}^{1-Y_0}]^{1-D_1} \\
& \left([R_{11}(X)^{Y_1} \{1 - R_{11}(X)\}^{1-Y_1}]^{Y_0} [R_{10}(X)^{Y_1} \{1 - R_{10}(X)\}^{1-Y_1}]^{1-Y_0} \right)^{D_1} \\
& \left([R_{01}(X)^{Y_1} \{1 - R_{01}(X)\}^{1-Y_1}]^{Y_0} [R_{00}(X)^{Y_1} \{1 - R_{00}(X)\}^{1-Y_1}]^{1-Y_0} \right)^{1-D_1}.
\end{aligned}$$

We will use γ to denote regular parametric submodels with γ_0 being the truth: e.g., $f(X; \gamma_0) = f(X)$. Here is our first lemma.

Lemma 5. *Suppose that assumption [E](#) holds. Then, the tangent space for $\bar{\theta}_L$ has the following form:*

$$\begin{aligned}
\mathcal{T} := & \left\{ \alpha_0(X) + \{D_1 - P(X)\} \alpha_1(X) + D_1 \{Y_0 - Q_1(X)\} \alpha_2(X) \right. \\
& + (1 - D_1) \{Y_0 - Q_0(X)\} \alpha_3(X) + D_1 Y_0 \{Y_1 - R_{11}(X)\} \alpha_4(X) \\
& \left. + D_1 (1 - Y_0) \{Y_1 - R_{10}(X)\} \alpha_5(X) + (1 - D_1) Y_0 \{Y_1 - R_{01}(X)\} \alpha_6(X) \right\}
\end{aligned}$$

$$+ (1 - D_1)(1 - Y_0)\{Y_1 - R_{00}(X)\}\alpha_7(X)\},$$

where α_j 's are all functions of X such that $\mathbb{E}\{\alpha_0(X)\} = 0$ and $\mathbb{E}\{\alpha_j^2(X)\} < \infty$ for $j = 0, 1, \dots, 7$.

Proof. The loglikelihood of regular parametric submodel is given by

$$\ell(\gamma) := \ell_0(\gamma) + \ell_{D_1}(\gamma) + \ell_{Y_0|D_1}(\gamma) + \ell_{Y_1|D_1, Y_0}(\gamma),$$

where

$$\ell_0(\gamma) := \log f(X; \gamma),$$

$$\ell_{D_1}(\gamma) := D_1 \log P(X; \gamma) + (1 - D_1) \log\{1 - P(X; \gamma)\},$$

$$\begin{aligned} \ell_{Y_0|D_1}(\gamma) = & D_1 \left[Y_0 \log Q_1(X; \gamma) + (1 - Y_0) \log\{1 - Q_1(X; \gamma)\} \right] \\ & + (1 - D_1) \left[Y_0 \log Q_0(X; \gamma) + (1 - Y_0) \log\{1 - Q_0(X; \gamma)\} \right], \end{aligned}$$

and

$$\begin{aligned} \ell_{Y_1|D_1, Y_0}(\gamma) := & D_1 \left(Y_0 [Y_1 \log R_{11}(X; \gamma) + (1 - Y_1) \log\{1 - R_{11}(X; \gamma)\}] \right. \\ & \left. + (1 - Y_0) [Y_1 \log R_{10}(X; \gamma) + (1 - Y_1) \log\{1 - R_{10}(X; \gamma)\}] \right) \\ & + (1 - D_1) \left(Y_0 [Y_1 \log R_{01}(X; \gamma) + (1 - Y_1) \log\{1 - R_{01}(X; \gamma)\}] \right. \\ & \left. + (1 - Y_0) [Y_1 \log R_{00}(X; \gamma) + (1 - Y_1) \log\{1 - R_{00}(X; \gamma)\}] \right). \end{aligned}$$

Therefore, the score at the truth has the following form:

$$\begin{aligned} S(Y_0, Y_1, D_1, X) := & \frac{1}{f(X)} \frac{\partial f(X; \gamma_0)}{\partial \gamma} + \frac{D_1 - P(X)}{P(X)\{1 - P(X)\}} \frac{\partial P(X; \gamma_0)}{\partial \gamma} \\ & + \frac{D_1\{Y_0 - Q_1(X)\}}{Q_1(X)\{1 - Q_1(X)\}} \frac{\partial Q_1(X; \gamma_0)}{\partial \gamma} + \frac{(1 - D_1)\{Y_0 - Q_0(X)\}}{Q_0(X)\{1 - Q_0(X)\}} \frac{\partial Q_0(X; \gamma_0)}{\partial \gamma} \\ & + \frac{D_1 Y_0\{Y_1 - R_{11}(X)\}}{R_{11}(X)\{1 - R_{11}(X)\}} \frac{\partial R_{11}(X; \gamma_0)}{\partial \gamma} + \frac{D_1(1 - Y_0)\{Y_1 - R_{10}(X)\}}{R_{10}(X)\{1 - R_{10}(X)\}} \frac{\partial R_{10}(X; \gamma_0)}{\partial \gamma} \\ & + \frac{(1 - D_1) Y_0\{Y_1 - R_{01}(X)\}}{R_{01}(X)\{1 - R_{01}(X)\}} \frac{\partial R_{01}(X; \gamma_0)}{\partial \gamma} + \frac{(1 - D_1)(1 - Y_0)\{Y_1 - R_{00}(X)\}}{R_{00}(X)\{1 - R_{00}(X)\}} \frac{\partial R_{00}(X; \gamma_0)}{\partial \gamma}, \end{aligned} \tag{42}$$

from which the lemma follows, because the derivatives are not restricted except for square integrability. \square

APPENDIX S-5. DERIVATION OF THE PATHWISE DERIVATIVES

Let

$$\begin{aligned}\bar{\theta}_{L,num} &:= \mathbb{E} \left[\{ \Pi_1(1, X) - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X) \} P(X) \right], \\ \bar{\theta}_{L,den} &:= \mathbb{E} \left[\{ 1 - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X) \} P(X) \right]\end{aligned}$$

so that $\bar{\theta}_L = \bar{\theta}_{L,num} / \bar{\theta}_{L,den}$. Define

$$\begin{aligned}F_{num}(Y_0, Y_1, D_1, X) &:= \{ \Pi_1(1, X) - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X) \} D_1 - \bar{\theta}_{L,num} \\ &\quad + D_1 [\{ Y_1 - \Pi_1(1, X) \} - \{ Y_0 - \Pi_0(1, X) \}] \\ &\quad - \frac{P(X)}{1 - P(X)} (1 - D_1) [\{ Y_1 - \Pi_1(0, X) \} - \{ Y_0 - \Pi_0(0, X) \}].\end{aligned}$$

Similarly, define

$$\begin{aligned}F_{den}(Y_0, Y_1, D_1, X) &:= \{ 1 - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X) \} D_1 - \bar{\theta}_{L,den} \\ &\quad - D_1 \{ Y_0 - \Pi_0(1, X) \} - \frac{P(X)}{1 - P(X)} (1 - D_1) [\{ Y_1 - \Pi_1(0, X) \} - \{ Y_0 - \Pi_0(0, X) \}].\end{aligned}$$

We will derive the pathwise derivatives of the numerator and denominator of $\bar{\theta}_L$ in a few steps. The following two lemmas show that $F_{num}(Y_0, Y_1, D_1, X)$ and $F_{den}(Y_0, Y_1, D_1, X)$ are the pathwise derivatives of $\bar{\theta}_{L,num}$ and $\bar{\theta}_{L,den}$, respectively.

Lemma 6. *Suppose assumption E is satisfied. Then, the pathwise derivative of $\bar{\theta}_{L,num}$ is given by $F_{num}(Y_0, Y_1, D_1, X)$.*

Proof. Using the fact that

$$\Pi_0(d, X) = Q_d(X) \quad \text{and} \quad \Pi_1(d, X) = Q_d(X)R_{d1}(X) + \{1 - Q_d(X)\}R_{d0}(X), \quad (43)$$

we can write

$$\begin{aligned}\bar{\theta}_{L,num} = \int & [Q_1(x)R_{11}(x) + \{1 - Q_1(x)\}R_{10}(x) - Q_1(x) \\ & - Q_0(x)R_{01}(x) - \{1 - Q_0(x)\}R_{00}(x) + Q_0(x)]P(x)f(x)dx.\end{aligned}$$

Therefore, the pathwise perturbation $\bar{\theta}_{L,num}(\gamma)$ of $\bar{\theta}_{L,num}$ is given by

$$\begin{aligned}\bar{\theta}_{L,num}(\gamma) = \int & [Q_1(x, \gamma)R_{11}(x, \gamma) + \{1 - Q_1(x, \gamma)\}R_{10}(x, \gamma) - Q_1(x, \gamma) \\ & - Q_0(x, \gamma)R_{01}(x, \gamma) - \{1 - Q_0(x, \gamma)\}R_{00}(x, \gamma) + Q_0(x, \gamma)]P(x, \gamma)f(x, \gamma)dx.\end{aligned}\quad (44)$$

Now, by straightforward algebra, $\partial\bar{\theta}_{L,num}(\gamma_0)/\partial\gamma$ is equal to the sum of the following terms (with each line corresponding to one term):

$$\begin{aligned}& \mathbb{E}\left[\{\Pi_1(1, X) - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}P(X)\frac{\partial f(X; \gamma_0)}{\partial\gamma}\frac{1}{f(X)}\right] \\ & \mathbb{E}\left[\{\Pi_1(1, X) - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}\frac{\partial P(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[\{R_{11}(X) - R_{10}(X) - 1\}P(X)\frac{\partial Q_1(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[\{R_{00}(X) - R_{01}(X) + 1\}P(X)\frac{\partial Q_0(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[Q_1(X)P(X)\frac{\partial R_{11}(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[\{1 - Q_1(X)\}P(X)\frac{\partial R_{10}(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[-Q_0(X)P(X)\frac{\partial R_{01}(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[-\{1 - Q_0(X)\}P(X)\frac{\partial R_{00}(X; \gamma_0)}{\partial\gamma}\right]\end{aligned}$$

We are looking for some $F(Y_0, Y_1, D_1, X)$ with mean zero that satisfies

$$\frac{\partial\bar{\theta}_{L,num}(\gamma_0)}{\partial\gamma} = \mathbb{E}\{F(Y_0, Y_1, D_1, X)S(Y_0, Y_1, D_1, X)\},\quad (45)$$

where $S(Y_0, Y_1, D_1, X)$ is the score described in equation (42). Using the fact that the variance of a binary variable is the “success” probability times that of “failure,” we know by

inspection that such $F(Y_0, Y_1, D_1, X)$ must be given by

$$\begin{aligned}
F(Y_0, Y_1, D_1, X) &:= \{\Pi_1(1, X) - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}P(X) - \bar{\theta}_{L,num} \\
&+ \{\Pi_1(1, X) - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}\{D_1 - P(X)\} \\
&+ \{R_{11}(X) - R_{10}(X) - 1\}P(X)\frac{D_1\{Y_0 - Q_1(X)\}}{P(X)} \\
&+ \{R_{00}(X) - R_{01}(X) + 1\}P(X)\frac{(1 - D_1)\{Y_0 - Q_0(X)\}}{1 - P(X)} \\
&+ Q_1(X)P(X)\frac{D_1Y_0\{Y_1 - R_{11}(X)\}}{P(X)Q_1(X)} \\
&+ \{1 - Q_1(X)\}P(X)\frac{D_1(1 - Y_0)\{Y_1 - R_{10}(X)\}}{P(X)\{1 - Q_1(X)\}} \\
&- Q_0(X)P(X)\frac{(1 - D_1)Y_0\{Y_1 - R_{01}(X)\}}{\{1 - P(X)\}Q_0(X)} \\
&- \{1 - Q_0(X)\}P(X)\frac{(1 - D_1)(1 - Y_0)\{Y_1 - R_{00}(X)\}}{\{1 - P(X)\}\{1 - Q_0(X)\}}.
\end{aligned}$$

However, simplifying $F(Y_0, Y_1, D_1, X)$ by using equation (43) yields

$$F(Y_0, Y_1, D_1, X) = F_{num}(Y_0, Y_1, D_1, X). \quad \square$$

Lemma 7. *Suppose assumption E is satisfied. Then, the pathwise derivative of $\bar{\theta}_{L,den}$ is given by $F_{den}(Y_0, Y_1, D_1, X)$.*

Proof. Using (43), we can write

$$\bar{\theta}_{L,den} = \int \{1 - Q_1(x) - Q_0(x)R_{01}(x) - \{1 - Q_0(x)\}R_{00}(x) + Q_0(x)\}P(x)f(x)dx.$$

Therefore, the pathwise perturbation of $\bar{\theta}_{L,den}$ is given by

$$\begin{aligned}
\bar{\theta}_{L,den}(\gamma) &:= \int \{1 - Q_1(x, \gamma) - Q_0(x, \gamma)R_{01}(x, \gamma) \\
&- \{1 - Q_0(x, \gamma)\}R_{00}(x, \gamma) + Q_0(x, \gamma)\}P(x, \gamma)f(x, \gamma)dx. \quad (46)
\end{aligned}$$

Hence, by simple algebra, $\partial\bar{\theta}_{L,den}(\gamma_0)/\partial\gamma$ is equal to the sum of the following terms (with each line corresponding to one term):

$$\begin{aligned} & \mathbb{E}\left[\{1 - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}P(X)\frac{\partial f(X; \gamma_0)}{\partial\gamma}\frac{1}{f(X)}\right] \\ & \mathbb{E}\left[\{1 - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}\frac{\partial P(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[-P(X)\frac{\partial Q_1(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[\{R_{00}(X) - R_{01}(X) + 1\}P(X)\frac{\partial Q_0(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[-Q_0(X)P(X)\frac{\partial R_{01}(X; \gamma_0)}{\partial\gamma}\right] \\ & \mathbb{E}\left[-\{1 - Q_0(X)\}P(X)\frac{\partial R_{00}(X; \gamma_0)}{\partial\gamma}\right]. \end{aligned}$$

Now, we are looking for some $F(Y_0, Y_1, D_1, X)$ with mean zero that satisfies

$$\frac{\partial\bar{\theta}_{L,den}(\gamma_0)}{\partial\gamma} = \mathbb{E}\{F(Y_0, Y_1, D_1, X)S(Y_0, Y_1, D_1, X)\}, \quad (47)$$

where $S(Y_0, Y_1, D_1, X)$ is the score described in equation (42). Using the fact that the variance of a binary variable is the “success” probability times that of “failure,” we know by inspection that such $F(Y_0, Y_1, D_1, X)$ must be given by

$$\begin{aligned} F(Y_0, Y_1, D_1, X) & := \{1 - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}P(X) - \bar{\theta}_{L,den} \\ & + \{1 - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}\{D_1 - P(X)\} \\ & - P(X)\frac{D_1\{Y_0 - Q_1(X)\}}{P(X)} \\ & + \{R_{00}(X) - R_{01}(X) + 1\}P(X)\frac{(1 - D_1)\{Y_0 - Q_0(X)\}}{1 - P(X)} \\ & - Q_0(X)P(X)\frac{(1 - D_1)Y_0\{Y_1 - R_{01}(X)\}}{\{1 - P(X)\}Q_0(X)} \\ & - \{1 - Q_0(X)\}P(X)\frac{(1 - D_1)(1 - Y_0)\{Y_1 - R_{00}(X)\}}{\{1 - P(X)\}\{1 - Q_0(X)\}}. \end{aligned}$$

However, simplifying $F(Y_0, Y_1, D_1, X)$ by using equation (43) yields

$$F(Y_0, Y_1, D_1, X) = F_{den}(Y_0, Y_1, D_1, X). \quad \square$$

In Section S-4, we have derived the scores at γ_0 and the tangent space. Also, we have shown that the pathwise derivatives of $\bar{\theta}_{L,num}$ and $\bar{\theta}_{L,den}$ are given by $F_{num}(Y_0, Y_1, D_1, X)$ and $F_{den}(Y_0, Y_1, D_1, X)$, respectively. By using these results, we now derive the pathwise derivative of $\bar{\theta}_L = \bar{\theta}_{L,num}/\bar{\theta}_{L,den}$ below.

Lemma 8. *Suppose assumption E is satisfied. Then, the pathwise derivative of $\bar{\theta}_L$ is given by*

$$G(Y_0, Y_1, D_1, X) := \frac{1}{\bar{\theta}_{L,den}} \left(F_{num}(Y_0, Y_1, D_1, X) - \bar{\theta}_L F_{den}(Y_0, Y_1, D_1, X) \right). \quad (48)$$

Proof. Since $\bar{\theta}_L = \bar{\theta}_{L,num}/\bar{\theta}_{L,den}$, we write the pathwise perturbation of $\bar{\theta}_L$ as $\bar{\theta}_L(\gamma) = \bar{\theta}_{L,num}(\gamma)/\bar{\theta}_{L,den}(\gamma)$, where the truth is denoted by γ_0 : see the section on the derivation of the pathwise derivatives of $\bar{\theta}_{L,num}$ and $\bar{\theta}_{L,den}$.

Now, we need to show that $G(Y_0, Y_1, D_1, X)$ satisfies

$$\frac{\partial \bar{\theta}_L(\gamma_0)}{\partial \gamma} = \mathbb{E} \{ G(Y_0, Y_1, D_1, X) S(Y_0, Y_1, D_1, X) \},$$

where $S(Y_0, Y_1, D_1, X)$ is the score described in equation (42) in Online Appendix S-4. But, noting that

$$\frac{\partial \bar{\theta}_L(\gamma_0)}{\partial \gamma} = \frac{1}{\bar{\theta}_{L,den}} \left(\frac{\partial \bar{\theta}_{L,num}(\gamma_0)}{\partial \gamma} - \bar{\theta}_L \frac{\partial \bar{\theta}_{L,den}(\gamma_0)}{\partial \gamma} \right), \quad (49)$$

it follows from the fact that the pathwise derivatives of $\bar{\theta}_{L,num}$ and $\bar{\theta}_{L,den}$ are given by $F_{num}(Y_0, Y_1, D_1, X)$ and $F_{den}(Y_0, Y_1, D_1, X)$, respectively: see Online Appendix S-5 on the derivations of the pathwise derivatives. \square

Lemma 9. *We have $F_{DID}(Y_0, Y_1, D_1, X) = G(Y_0, Y_1, D_1, X)$, where $F_{DID}(Y_0, Y_1, D_1, X)$ is defined in equation (17) in Theorem 4.*

Proof. This is a simple algebraic result. Specifically, if we plug $F_{num}(Y_0, Y_1, D_1, X)$ and $F_{den}(Y_0, Y_1, D_1, X)$ into equation (48), then it follows that

$$G(Y_0, Y_1, D_1, X) = \frac{1}{\bar{\theta}_{L,den}} \left(H_1(Y_0, Y_1, D_1, X) + H_2(Y_0, Y_1, D_1, X) + H_3(Y_0, Y_1, D_1, X) \right), \quad (50)$$

where

$$H_1(Y_0, Y_1, D_1, X) := D_1 [\{\Pi_1(1, X) - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\} \\ - \bar{\theta}_L \{1 - \Pi_0(1, X) - \Pi_1(0, X) + \Pi_0(0, X)\}],$$

$$H_2(Y_0, Y_1, D_1, X) := D_1 [\{Y_1 - \Pi_1(1, X)\} - (1 - \bar{\theta}_L) \{Y_0 - \Pi_0(1, X)\}]$$

$$H_3(Y_0, Y_1, D_1, X) := (\bar{\theta}_L - 1) \frac{P(X)}{1 - P(X)} (1 - D_1) [\{Y_1 - \Pi_1(0, X)\} - \{Y_0 - \Pi_0(0, X)\}].$$

Further simplifications yield the form of $F_{DID}(Y_0, Y_1, D_1, X)$ defined in equation (17) in Theorem 4. In addition, the form of $F_{DID}(Y_0, Y)1, D_1, X$ in (18) can be obtained by manipulating the terms in (17). \square

APPENDIX S-6. INFERENCE FOR THE EVENT-STUDY PERSUASION RATE

The asymptotic influence function for regular and asymptotically linear estimators of the multi-period parameters can be derived from our previous discussion on the two-period case. We will explain this by focusing on $\hat{\theta}_{ESPR,DR}(j)$, which we consider the most comprehensive summary parameter.

Let $\theta_{num}(s, s + j)$ and $\theta_{den}(s, s + j)$ be the numerator and the denominator of $\theta(s, s + j)$, respectively, as before. Let $\mathcal{D}_{s,s+j} := (Y_{s-1}, Y_{s+j}, S, X)$, and define

$$G_{num}(\mathcal{D}_{s,s+j}) := \mathbb{1}(S = s)(Y_{s+j} - Y_{s-1}) - \frac{\mathbb{P}(S = s \mid X, \bar{S}_s = 1)}{\mathbb{P}(S = \infty \mid X, \bar{S}_s = 1)} (Y_{s+j} - Y_{s-1}) \mathbb{1}(S = \infty),$$

$$G_{den}(\mathcal{D}_{s,s+j}) := \mathbb{1}(S = s)(1 - Y_{s-1}) - \frac{\mathbb{P}(S = s \mid X, \bar{S}_s = 1)}{\mathbb{P}(S = \infty \mid X, \bar{S}_s = 1)} (Y_{s+j} - Y_{s-1}) \mathbb{1}(S = \infty),$$

$$G_{adj}(\mathcal{D}_{s,s+j}) := - \left\{ \mathbb{1}(S = s) - \mathbb{1}(S = \infty) \frac{\mathbb{P}(S = s \mid X, \bar{S}_s = 1)}{\mathbb{P}(S = \infty \mid X, \bar{S}_s = 1)} \right\} \Delta_{s,s+j}(\infty, X),$$

where

$$\Delta_{s,s+j}(\infty, X) := \mathbb{P}(Y_{s+j} = 1 \mid S = \infty, X) - \mathbb{P}(Y_{s-1} = 1 \mid S = \infty, X).$$

We then have the following theorem, for which we consider a random sample of size n , $\{(Y_{i0}, Y_{i1}, \dots, Y_{iT}, S_i, X_i) : i = 1, \dots, n\}$. Let $p_s := \mathbb{P}(S_i = s)$, and $\hat{p}_s := n^{-1} \sum_{i=1}^n \mathbb{1}(S_i = s)$. Further, for $r \in \{\text{num}, \text{den}\}$, define

$$\begin{aligned} \tilde{P} &:= \left[\frac{p_1}{\mathbb{P}(S_i=1 \text{ or } S_i=\infty)} \quad \dots \quad \frac{p_{T-j}}{\mathbb{P}(S_i=T-j \text{ or } S_i=\infty)} \right]^\top, \\ \Theta_r &:= \left[\theta_r(1, 1+j) \quad \dots \quad \theta_r(T-j, T) \right]^\top. \end{aligned}$$

Finally, let

$$\mathbf{Q}_i := \left[\mathbf{G}_{\text{num},i}^\top \quad \mathbf{G}_{\text{den},i}^\top \quad \mathbf{H}_i^\top \right]^\top,$$

where for $r \in \{\text{num}, \text{den}\}$,

$$\begin{aligned} \mathbf{G}_{r,i} &:= \left[G_r(\mathcal{D}_{i,1,1+j}) + G_{\text{adj}}(\mathcal{D}_{i,1,1+j}) \quad \dots \quad G_r(\mathcal{D}_{i,T-j,T}) + G_{\text{adj}}(\mathcal{D}_{i,T-j,T}) \right]^\top \\ \mathbf{H}_i &:= \left[\mathbb{1}(S_i = 1) - p_1 \quad \dots \quad \mathbb{1}(S_i = T-j) - p_{T-j} \right]^\top. \end{aligned}$$

Theorem 7. For $r \in \{\text{num}, \text{den}\}$, let $\hat{\theta}_r(s, s+j)$ be an estimator of $\theta_r(s, s+j)$ based on a subsample $\{(Y_{is-1}, Y_{is+j}, \mathbb{1}(S_i = s), X_i) \mathbb{1}(S_i = 1 \text{ or } S_i = \infty) : i = 1, 2, \dots, n\}$. If they are regular and asymptotically linear (conditional on $\mathbb{1}(S_i = s) + \mathbb{1}(S_i = \infty) = 1$), then

$$\sqrt{n} \left(\frac{\sum_{s=1}^{T-j} \hat{p}_s \hat{\theta}_{\text{num}}(s, s+j)}{\sum_{s=1}^{T-j} \hat{p}_s \hat{\theta}_{\text{den}}(s, s+j)} - \frac{\sum_{s=1}^{T-j} p_s \theta_{\text{num}}(s, s+j)}{\sum_{s=1}^{T-j} p_s \theta_{\text{den}}(s, s+j)} \right) \xrightarrow{d} N\left(0, J P \Sigma P^\top J^\top\right),$$

where

$$\begin{aligned} J &:= \left[\frac{1}{\sum_{s=1}^{T-j} p_s \theta_{\text{den}}(s, s+j)} \quad - \frac{\sum_{s=1}^{T-j} p_s \theta_{\text{den}}(s, s+j)}{\{\sum_{s=1}^{T-j} p_s \theta_{\text{den}}(s, s+j)\}^2} \right], \\ P &:= \begin{bmatrix} \tilde{P}^\top & \mathbf{O}^\top & \Theta_{\text{num}}^\top \\ \mathbf{O}^\top & \tilde{P}^\top & \Theta_{\text{den}}^\top \end{bmatrix} \\ \Sigma &:= \mathbb{E}(\mathbf{Q}_1 \mathbf{Q}_1^\top), \end{aligned}$$

where \mathbf{O} is a $(T-j) \times 1$ vector of zeroes.

Therefore, the DR estimator $\hat{\theta}_{ESPR,DR}(j)$ will be asymptotically normal with the asymptotic variance provided in theorem 7, as long as its component estimators are regular and asymptotically linear.

Proof of Theorem 7: By similar calculations to theorem 4 but by using the conditional likelihood given $\mathbb{1}(S_i = s \text{ or } S_i = \infty) = 1$, we know that a regular and asymptotically linear estimator $\hat{\theta}_r(s, s + j)$ must have the following asymptotic expansion:

$$\begin{aligned} \sqrt{n}\{\hat{\theta}_r(s, s + j) - \theta_r(s, s + j)\} \\ = \frac{1}{\mathbb{P}(S_i = s \text{ or } S_i = \infty)} \frac{1}{\sqrt{n}} \sum_{i=1}^n \{G_r(\mathcal{D}_{i,s,s+j}) + G_{adj}(\mathcal{D}_{i,s,s+j})\} + o_p(1), \end{aligned}$$

where $\mathcal{D}_{i,s,s+j} := (Y_{is-1}, Y_{is+j}, \mathbb{1}(S_i = s), X_i)$. Therefore, let $p_s := \mathbb{P}(S_i = s)$, and we have

$$\sqrt{n} \sum_{s=1}^{T-j} p_s \{\hat{\theta}_r(s, s + j) - \theta_r(s, s + j)\} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{P}^\top \mathbf{G}_{r,i} + o_p(1).$$

Note also that

$$\sqrt{n} \sum_{s=1}^{T-j} \theta_r(s, s + j) (\hat{p}_s - p_s) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \Theta_r^\top \mathbf{H}_i.$$

Therefore, it follows that

$$\begin{aligned} \sqrt{n} \left\{ \sum_{s=1}^{T-j} \hat{p}_s \hat{\theta}_r(s, s + j) - \sum_{s=1}^{T-j} p_s \theta_r(s, s + j) \right\} \\ = \sqrt{n} \sum_{s=1}^{T-j} [p_s \{\hat{\theta}_r(s, s + j) - \theta_r(s, s + j)\} + \theta_r(s, s + j) (\hat{p}_s - p_s)] + O_p(n^{-1/2}) \\ = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\tilde{P}^\top \mathbf{G}_{r,i} + \Theta_r^\top \mathbf{H}_i) + o_p(1). \end{aligned} \tag{51}$$

Now, use the expressions in equation (51) for $r \in \{num, den\}$ together with the delta method. □